

Санкт-Петербургский Государственный Университет
Математико-механический факультет

Математическое обеспечение и администрирование информационных
систем

Шкуратов Илья Андреевич

Оптимизация запросов с обобщёнными табличными выражениями

Магистерская диссертация

Научный руководитель:
д. ф.-м. н., профессор Новиков Б. А.

Рецензент:
Яцышин И. А.

Санкт-Петербург
2017

SAINT-PETERSBURG STATE UNIVERSITY

Software and Administration of Information Systems

Ilia Shkuratov

Optimization of queries with common table expressions

Graduation Thesis

Scientific supervisor:
professor Boris Novikov

Reviewer:
Ilya Yatsishin

Saint-Petersburg
2017

Оглавление

Введение	4
1. Постановка задачи	5
2. Оптимизация обобщённых табличных выражений	6
3. Анализ планов запросов с ОТВ в PostgreSQL	8
3.1. Потеря порядка сортировки	8
3.2. «Проталкивание» предикатов	10
4. Оптимизатор/планировщик PostgreSQL	17
5. Возможность реализации алгоритма встраивания ОТВ в PostgreSQL	19
5.1. Сравнение планов со встроенными ОТВ	21
6. Мнение сообщества разработчиков PostgreSQL о встраивании ОТВ	26
Заключение	27
Список литературы	28

Введение

Обобщённые табличные выражения (ОТВ) были введены в стандарте SQL/Foundation ISO/IEC 9075-2:1999. Они увеличивают выразительность языка SQL, позволяя писать рекурсивные запросы, а также представляют средства для лучшего структурирования запроса, путём разбиения его на небольшие именованные подзапросы.

На обобщённые табличные выражения можно ссылаться по имени как из других ОТВ, которые объявлены после, так и из тела запроса. В запросе может быть несколько ссылок на ОТВ, аналогично ссылкам на таблицу. При этом, семантика ОТВ накладывает ограничение: результат, получаемый по этой ссылке должен быть таким, как-будто ОТВ выполнилось один раз.

В данной работе рассматриваются нерекурсивные ОТВ и возможность их оптимизации в PostgreSQL. Оптимизатор PostgreSQL реализует семантику простым способом — он материализует результат выполнения ОТВ, а затем использует его при каждом обращении. При этом он оптимизирует подзапрос соответствующий ОТВ в изоляции от остальной части, что может приводить к неэффективным планам в некоторых сценариях.

1. Постановка задачи

В связи с особенностями реализации обобщённых табличных выражений (ОТВ) в PostgreSQL, целью данной работы были поставлены следующие задачи:

- исследовать механизм работы оптимизатора с ОТВ в текущей версии PostgreSQL,
- изучить существующие методы оптимизации ОТВ,
- предложить способ их реализации в текущей инфраструктуре оптимизатора.

2. Оптимизация обобщённых табличных выражений

Популярные СУБД, такие как DB2, Oracle Database, SQL Server поддерживают для ОТВ тот же набор оптимизаций, что и для обычных подзапросов или представлений ¹. В случае если оптимизатор считает, что замена ссылки на ОТВ соответствующим подзапросом (встраивание) не изменит семантику запроса, он рассматривает планы с данной альтернативой. Встраивание позволяет использовать все методы оптимизации, обычно применяемые к подзапросам [7, 2, 1, 4, 5]:

- проталкивание предикатов,
- поднятие предикатов,
- слияние с внешним запросом (unnesting),
- слияние с другим подзапросом (coalescing),
- замена на оконную функцию.

Однако встраивание не возможно, когда ОТВ содержит недетерминированные (volatile) функции или функции с доступом ко внешним ресурсам, так как это может нарушить семантику. Кроме того, не смотря на широкий спектр возможных оптимизаций, встраивание не всегда будет оказываться выгодным [8, 3].

Чтобы принять оптимальное решение, нам нужно рассмотреть все возможные планы получающиеся при встраивании ОТВ вместо ссылки [8]. Однако узнать стоимость таких планов мы сможем только на уровне наименьшего общего предка (НОК) ссылок на ОТВ [3]. Это так в силу того, что мы не можем сразу добавить стоимость выполнения и материализации исходного ОТВ (C^{CTE}) к планам, которые используют обращение по ссылке, так как эта стоимость будет «разделяться» между всеми такими обращениями. А в случае плана, которые использует только встроенные ОТВ, мы не можем знать нужно ли добавлять

¹<http://modern-sql.com/feature/with/performance>

стоимость C^{CTE} , так как в другой части запроса может быть выгодно использовать обращение к материализованному ОТВ.

Дойдя до НОК, мы сможем учесть C^{CTE} и сравнить стоимости планов, содержащих встроенные ОТВ, с остальными планами и отбросить те, что имеют большую стоимость.

Понятно, что количество дополнительных планов, которые нам нужно будет рассмотреть, экспоненциально зависит от количества ссылок на ОТВ, и при их большом количестве, нужно использовать эвристические подходы.

Кроме описанной выше оптимизации, разработчики Orca [9] реализовали в своей системе проталкивание предикатов в ОТВ и учёт физических свойств данных в местах использования [8]. При проталкивании, предикаты из разных контекстов (относящиеся к разным ссылкам на ОТВ) объединяются дизъюнкцией и добавляется к подзапросу ОТВ. При дальнейшей оптимизации ОТВ они, по возможности, будут проталкиваться дальше. При этом оригинальные предикаты остаются на месте, чтобы сохранить семантику запроса. Учёт физических свойств позволяет протолкнуть общие требования встречающиеся в контексте ссылок на ОТВ, в само ОТВ. Например, если во всех контекстах требуется определённый порядок записей или партиционирование по определённому ключу, то соответствующие операции будут применены один раз, при выполнении ОТВ.

3. Анализ планов запросов с ОТВ в PostgreSQL

Для проведения анализа использовалась демонстрационная база данных Postgres Professional [12] маленького размера и последняя версия PostgreSQL [10]. Эксперименты проводились на машине со следующими характеристиками:

RAM	8 Gb
CPU	Intel® Core™ i5-4200U CPU @ 1.60GHz × 4
OS	Ubuntu 16.04 LTS, 64-bit
HDD	ST500LM021-1KJ152

Таблица 1: Конфигурация экспериментальной машины

Тематикой демонстрационной базы являются авиаперевозки. Схема базы данных показана на Рис. 1. Мы не будем описывать данные, так как схемы должно быть достаточно для понимания представленных запросов. При необходимости, можно получить подробное описание на сайте с базой [12].

Рассмотрим несколько примеров, которые иллюстрируют недостатки отказа оптимизации ОТВ совместно с остальной частью запроса.

3.1. Потеря порядка сортировки

Посмотрим на Запрос 1, который выводит количество мест в самолёте, сгруппированное по классу обслуживания. Сравним его план выполнения (План 1), с планом выполнения (План 2) семантически эквивалентного Запроса 2, использующего обобщённое табличное выражение.

Обратим внимание на повторную сортировку результата в Плане 2 (строка 10). Ввиду того, что оптимизатор не запоминает порядок записей в ОТВ, он вынужден добавлять в план сортировку, даже когда в этом нет необходимости. При работе с подзапросами, эта информация сохраняется, поэтому оптимизатор может сгенерировать план без лишней сортировки.


```

1 SELECT  s2.aircraft_code,
2         string_agg (s2.fare_conditions || '(' || s2.num::text || ')',
3                   ',') as fare_conditions
4 FROM    (
5         SELECT  s.aircraft_code, s.fare_conditions, count(*) as num
6         FROM    seats s
7         GROUP BY s.aircraft_code, s.fare_conditions
8         ORDER BY s.aircraft_code, s.fare_conditions
9     ) s2
10 GROUP BY s2.aircraft_code
11 ORDER BY s2.aircraft_code;

```

Запрос 1: Количество мест в самолёте, с учётом класса обслуживания.

```

1 GroupAggregate (cost=32.34..33.49 rows=27 width=36)
2   Group Key: s.aircraft_code
3   → Sort (cost=32.34..32.41 rows=27 width=20)
4     Sort Key: s.aircraft_code, s.fare_conditions
5     → HashAggregate (cost=31.43..31.70 rows=27 width=20)
6       Group Key: s.aircraft_code, s.fare_conditions
7       → Seq Scan on seats s (cost=0.00..21.39 rows=1339 width=12)

```

План 1: План к Запросу 1.

```

1 WITH s2 AS (
2   SELECT  s.aircraft_code, s.fare_conditions, count(*) as num
3   FROM    seats s
4   GROUP BY s.aircraft_code, s.fare_conditions
5   ORDER BY s.aircraft_code, s.fare_conditions
6 )
7 SELECT  s2.aircraft_code,
8         string_agg (s2.fare_conditions || '(' || s2.num::text || ')',
9                   ',') as fare_conditions
10 FROM    s2
11 GROUP BY s2.aircraft_code
12 ORDER BY s2.aircraft_code;

```

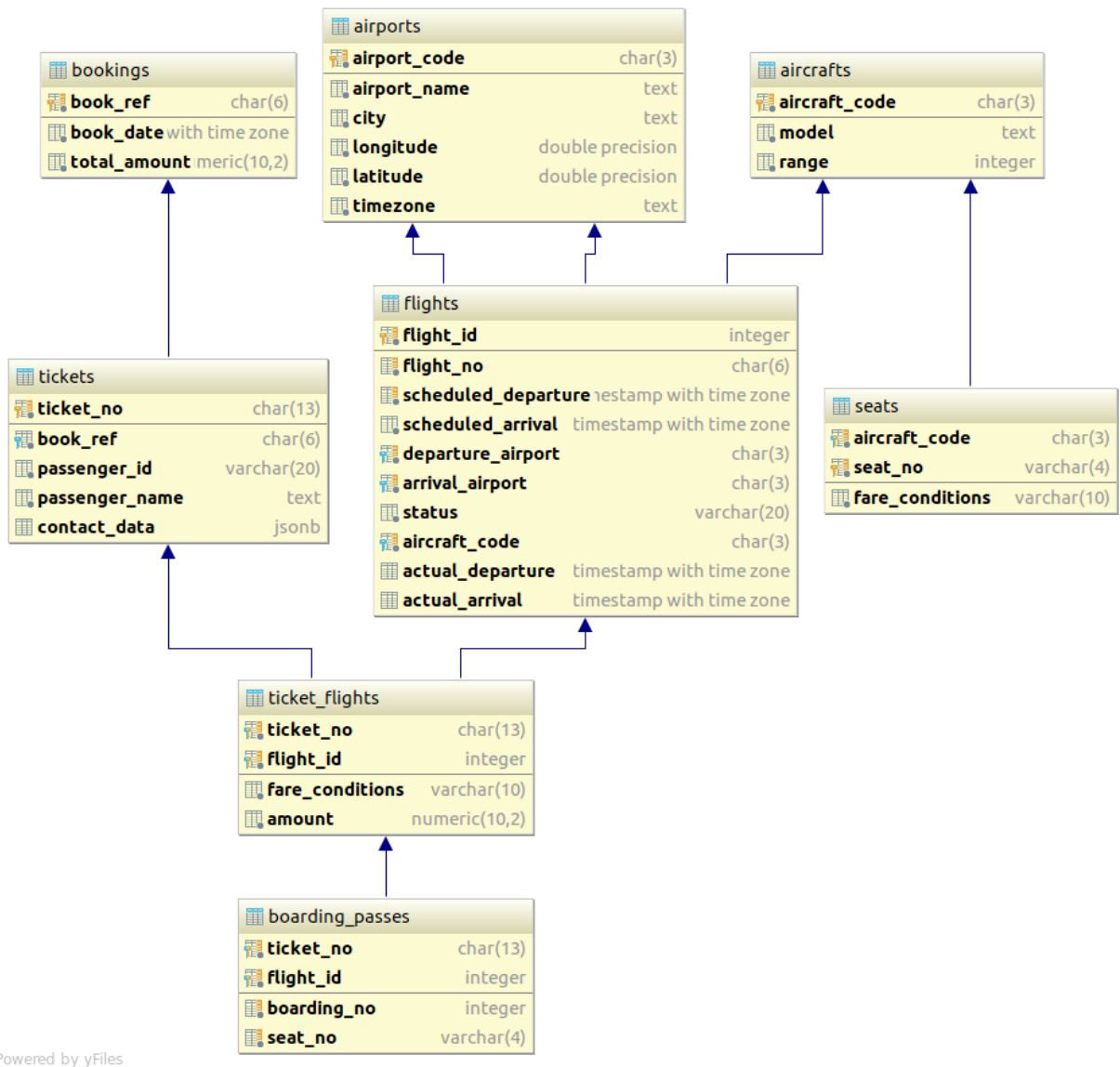
Запрос 2: Запрос с ОТВ, эквивалентный по результату Запросу 1.

```

1 GroupAggregate (cost=33.59..34.47 rows=27 width=48)
2   Group Key: s2.aircraft_code
3   CTE s2
4     → Sort (cost=32.34..32.41 rows=27 width=20)
5       Sort Key: s.aircraft_code, s.fare_conditions
6       → HashAggregate (cost=31.43..31.70 rows=27 width=20)
9         Group Key: s.aircraft_code, s.fare_conditions
10        → Seq Scan on seats s (cost=0.00..21.39 rows=1339 width=12)
11     → Sort (cost=1.18..1.25 rows=27 width=62)
12       Sort Key: s2.aircraft_code
13       → CTE Scan on s2 (cost=0.00..0.54 rows=27 width=62)

```

План 2: План к Запросу 2.



Powered by yFiles

Рис. 1: Схема тестовой базы данных.

3.2. «Проталкивание» предикатов

Рассмотрим запрос, который находит все стыковочные рейсы из Домодедово в Пулково без ограничений по времени вылета и продолжительности пересадки. В этот раз мы начнём с запроса, использующего ОТВ. Так нам будет легче понять его структуру и проследить влияние оптимизации на план, при поэтапном встраивании ОТВ.

В исходном Запросе 3 используются три ОТВ:

- **FlightSeatsCount** подсчитывает количество мест на каждом судне.
- **FlightSeatsOccupied** подсчитывает количество занятых мест на

```

1 WITH FlightSeatsCount as (
2   SELECT aircraft_code, fare_conditions, count(seat_no) AS total
3   FROM seats
4   GROUP BY (aircraft_code, fare_conditions)
5 ), FlightSeatsOccupied as (
6   SELECT
7     fl.flight_id, fl.departure_airport, fl.arrival_airport,
8     fl.scheduled_departure, fl.scheduled_arrival, fl.aircraft_code,
9     tf.fare_conditions, count(ticket_no) AS occupied
10  FROM flights fl
11  JOIN ticket_flights tf ON fl.flight_id = tf.flight_id
12  GROUP BY (fl.flight_id, tf.fare_conditions)
13 ), FlightSeatsAvailable as (
14  SELECT fo.flight_id, fo.departure_airport, fo.arrival_airport,
15         fo.scheduled_departure, fo.scheduled_arrival,
16         fc.fare_conditions, (fc.total - fo.occupied) AS available
17  FROM FlightSeatsCount fc
18  JOIN FlightSeatsOccupied fo
19       ON fc.aircraft_code = fo.aircraft_code AND
20         fc.fare_conditions = fo.fare_conditions
21  WHERE fc.total - fo.occupied > 0
22 )
23 SELECT
24   flo.flight_id, flo.departure_airport, flo.arrival_airport,
25   flo.available, flo.fare_conditions,
26   flo.scheduled_departure, flo.scheduled_arrival,
27   fli.flight_id, fli.departure_airport, fli.arrival_airport,
28   fli.available, fli.fare_conditions,
29   fli.scheduled_departure, fli.scheduled_arrival
30 FROM FlightSeatsAvailable flo
31 JOIN FlightSeatsAvailable fli ON
32     flo.arrival_airport = fli.departure_airport AND
33     flo.scheduled_arrival < fli.scheduled_departure
34 WHERE flo.departure_airport = 'DME' AND fli.arrival_airport = 'LED'

```

Запрос 3: Нахождение стыковочных рейсов из Домодедово в Пулково.

полёте.

- **FlightSeatsAvailable** использует результаты предыдущих вычислений, чтобы посчитать количество свободных мест на полёте.

Выполнение данного запроса занимает значительной время — около двух минут. Большая часть времени уходит на самое внешнее соединение, при котором отфильтровывается 42465913 строк. Также, мы можем заметить, что при выполнении **FlightSeatsOccupied** происходит внешняя сортировка. Мы можем сильно сократить количество обрабатываемых данных, если сразу отфильтруем ненужные полёты.

Встраивание **FlightSeatsAvailable** и **FlightSeatsOccupied** (Запрос 4)

```

1 Nested Loop (cost=296683.39..296685.43 rows=1 width=196) (actual time
  =5956.278..232121.631 rows=112448 loops=1)
2   Join Filter: ((flo.scheduled_arrival < fli.scheduled_departure) AND (flo.
  arrival_airport = fli.departure_airport))
3   Rows Removed by Join Filter: 42465913
4   CTE flightseatscount
5     -> HashAggregate (cost=31.43..31.70 rows=27 width=20)
6         Group Key: seats.aircraft_code, seats.fare_conditions
7         -> Seq Scan on seats (cost=0.00..21.39 rows=1339 width=15)
8   CTE flightseatsoccupied
9     -> Finalize GroupAggregate (cost=227984.02..291146.85 rows=200064 width
  =48) (actual time=3975.011..5183.306 rows=74697 loops=1)
10        Group Key: fl.flight_id, tf.fare_conditions
11        -> Gather Merge (cost=227984.02..286145.25 rows=400128 width=48)
12            Workers Planned: 2
13            -> Partial GroupAggregate (cost=226983.99..238960.52 rows
  =200064 width=48)
14                Group Key: fl.flight_id, tf.fare_conditions
15                -> Sort (cost=226983.99..229477.97 rows=997589 width
  =54)
16                    Sort Key: fl.flight_id, tf.fare_conditions
17                    Sort Method: external merge  Disk: 62472kB
18                -> Hash Join (cost=2764.48..59382.80 rows=997589
  width=54)
19                    Hash Cond: (tf.flight_id = fl.flight_id)
20                    -> Parallel Seq Scan on ticket_flights tf
  (cost=0.00..29927.89 rows=997589 width
  =26)
21                    -> Hash (cost=1474.88..1474.88 rows=66688
  width=32)
22                        -> Seq Scan on flights fl (cost
  =0.00..1474.88 rows=66688 width
  =32)
23   CTE flightseatsavailable
24     -> Hash Join (cost=0.95..5504.84 rows=45 width=98) (actual time
  =3976.147..5383.548 rows=71223 loops=1)
25         Hash Cond: ((fo.aircraft_code = fc.aircraft_code) AND ((fo.
  fare_conditions)::text = (fc.fare_conditions)::text))
26         Join Filter: ((fc.total - fo.occupied) > 0)
27         Rows Removed by Join Filter: 3474
28         -> CTE Scan on flightseatsoccupied fo (cost=0.00..4001.28 rows
  =200064 width=114) (actual time=3975.015..5284.268 rows=74697
  loops=1)
29         -> Hash (cost=0.54..0.54 rows=27 width=62)
30             -> CTE Scan on flightseatscount fc (cost=0.00..0.54 rows=27
  width=62)
31     -> CTE Scan on flightseatsavailable flo (cost=0.00..1.01 rows=1 width=98)
32         Filter: (departure_airport = 'DME'::bpchar)
33         Rows Removed by Filter: 62424
34     -> CTE Scan on flightseatsavailable fli (cost=0.00..1.01 rows=1 width=98)
35         Filter: (arrival_airport = 'LED'::bpchar)
36         Rows Removed by Filter: 66384

```

План 3: План выполнения Запроса 3.

может сократить число полётов для которых подсчитываются свободные места, так как позволяет оптимизатору протолкнуть предикаты к листьям плана (План 4), к сканированию таблицы `flights`.

После проведённого рефакторинга, появилось много дублирующегося кода, запрос стал значительно менее понятным, однако время его выполнения, благодаря «проталкиванию» предикатов, уменьшилось радикально — с 2 минут до 2 секунд. Как видно из данного примера, возможность проталкивание предикатов может кардинально сократить время выполнения запроса.

```

1 WITH FlightSeatsCount as (
2   SELECT aircraft_code, fare_conditions, count(seat_no) AS total
3   FROM seats
4   GROUP BY (aircraft_code, fare_conditions)
5 )
6 SELECT
7   flo.flight_id, flo.departure_airport, flo.arrival_airport,
8   flo.available, flo.fare_conditions,
9   flo.scheduled_departure, flo.scheduled_arrival,
10  fli.flight_id, fli.departure_airport, fli.arrival_airport,
11  fli.available, fli.fare_conditions,
12  fli.scheduled_departure, fli.scheduled_arrival
13 FROM (
14  SELECT fo.flight_id, fo.departure_airport, fo.arrival_airport,
15         fo.scheduled_departure, fo.scheduled_arrival,
16         fc.fare_conditions, (fc.total - fo.occupied) AS available
17  FROM FlightSeatsCount fc
18  JOIN (
19    SELECT
20      fl.flight_id, fl.departure_airport, fl.arrival_airport,
21      fl.scheduled_departure, fl.scheduled_arrival, fl.aircraft_code,
22      tf.fare_conditions, count(ticket_no) AS occupied
23    FROM flights fl
24    JOIN ticket_flights tf ON fl.flight_id = tf.flight_id
25    GROUP BY (fl.flight_id, tf.fare_conditions)
26  ) fo
27    ON fc.aircraft_code = fo.aircraft_code AND
28    fc.fare_conditions = fo.fare_conditions
29  WHERE fc.total - fo.occupied > 0
30 ) flo
31 JOIN (
32  SELECT fo.flight_id, fo.departure_airport, fo.arrival_airport,
33         fo.scheduled_departure, fo.scheduled_arrival,
34         fc.fare_conditions, (fc.total - fo.occupied) AS available
35  FROM FlightSeatsCount fc
36  JOIN (
37    SELECT
38      fl.flight_id, fl.departure_airport, fl.arrival_airport,
39      fl.scheduled_departure, fl.scheduled_arrival, fl.aircraft_code,
40      tf.fare_conditions, count(ticket_no) AS occupied
41    FROM flights fl
42    JOIN ticket_flights tf ON fl.flight_id = tf.flight_id
43    GROUP BY (fl.flight_id, tf.fare_conditions)
44  ) fo
45    ON fc.aircraft_code = fo.aircraft_code AND
46    fc.fare_conditions = fo.fare_conditions
47  WHERE fc.total - fo.occupied > 0
48 ) fli ON
49         flo.arrival_airport = fli.departure_airport AND
50         flo.scheduled_arrival < fli.scheduled_departure
51 WHERE flo.departure_airport = 'DME' AND fli.arrival_airport = 'LED';

```

Запрос 4: Нахождение стыковочных рейсов, со встроенными FlightSeatsAvailable и FlightSeatsOccupied.

```

1 Hash Join (cost=98883.72..99342.67 rows=1 width=148)
2   Hash Cond: ((fl.aircraft_code = fc.aircraft_code) AND ((tf.fare_conditions)
3     ::text = (fc.fare_conditions)::text))
4   Join Filter: ((fc.total - (count(tf.ticket_no))) > 0)
5   CTE flightseatscount
6     -> HashAggregate (cost=31.43..31.70 rows=27 width=20)
7       Group Key: seats.aircraft_code, seats.fare_conditions
8         -> Seq Scan on seats (cost=0.00..21.39 rows=1339 width=15)
9     -> Hash Join (cost=98851.07..99309.38 rows=83 width=130)
10       Hash Cond: (fl.arrival_airport = fl_1.departure_airport)
11       Join Filter: (fl.scheduled_arrival < fl_1.scheduled_departure)
12       -> Finalize HashAggregate (cost=50267.56..50459.02 rows=19146 width
13         =48)
14         Group Key: fl.flight_id, tf.fare_conditions
15         -> Gather (cost=45959.71..49980.37 rows=38292 width=48)
16           Workers Planned: 2
17           -> Partial HashAggregate (cost=44959.71..45151.17 rows
18             =19146 width=48)
19             Group Key: fl.flight_id, tf.fare_conditions
20             -> Hash Join (cost=1721.38..44243.69 rows=95469
21               width=54)
22               Hash Cond: (tf.flight_id = fl.flight_id)
23                 -> Parallel Seq Scan on ticket_flights tf (
24                   cost=0.00..29927.89 rows=997589 width=26)
25                 -> Hash (cost=1641.60..1641.60 rows=6382
26                   width=32)
27                   -> Seq Scan on flights fl (cost
28                     =0.00..1641.60 rows=6382 width=32)
29                     Filter: (departure_airport = 'DME
30                       '::bpchar)

```

План 4: План выполнения Запроса 4. Часть 1.

```

23     -> Hash (cost=48583.48..48583.48 rows=3 width=82)
24         -> Hash Join (cost=48263.67..48583.48 rows=3 width=82)
25             Hash Cond: ((fl_1.aircraft_code = fc_1.aircraft_code) AND
26                 ((tf_1.fare_conditions)::text = (fc_1.fare_conditions)
27                 ::text))
26             Join Filter: ((fc_1.total - (count(tf_1.ticket_no))) > 0)
27         -> Finalize HashAggregate (cost=48262.72..48378.97 rows
28             =11625 width=48)
29             Group Key: fl_1.flight_id, tf_1.fare_conditions
30                 -> Gather (cost=45647.10..48088.35 rows=23250
31                     width=48)
32                     Workers Planned: 2
33                     -> Partial HashAggregate (cost
34                         =44647.10..44763.35 rows=11625 width=48)
35                         Group Key: fl_1.flight_id, tf_1.
36                         fare_conditions
37                         -> Hash Join (cost=1690.04..44212.35
38                             rows=57966 width=54)
39                             Hash Cond: (tf_1.flight_id = fl_1.
40                                 flight_id)
41                             -> Parallel Seq Scan on
42                                 ticket_flights tf_1 (cost
43                                     =0.00..29927.89 rows=997589
44                                     width=26)
45                             -> Hash (cost=1641.60..1641.60
46                                 rows=3875 width=32)
47                                 -> Seq Scan on flights fl_1
48                                     (cost=0.00..1641.60
49                                     rows=3875 width=32)
50                                     Filter: (
51                                         arrival_airport =
52                                         'LED'::bpchar)
53         -> Hash (cost=0.54..0.54 rows=27 width=62)
54             -> CTE Scan on flightseatscount fc_1 (cost
55                 =0.00..0.54 rows=27 width=62)
56 -> Hash (cost=0.54..0.54 rows=27 width=62)
57     -> CTE Scan on flightseatscount fc (cost=0.00..0.54 rows=27 width
58         =62)

```

План 4: План выполнения Запроса 4. Часть 2.

4. Оптимизатор/планировщик PostgreSQL

В PostgreSQL нет чёткого разделения на оптимизатор и планировщик. По сути, это один компонент, который переписывает запрос и тут же строит план его выполнения. В дальнейшем, мы будем называть этот компонент оптимизатором.

На вход оптимизатору приходит дерево запроса, прошедшее семантический анализ и представляющее собой логический план выполнения запроса. На выходе получается физический план, который также является деревом.

Процесс планирования выполняется одним потоком. Оптимизатор рекурсивно обрабатывает подзапросы, строя план «снизу-вверх». Упрощённая последовательность действий, выполняемая оптимизатором, представлена в Таблице 2. Также в ней указаны имена некоторых важных функций, на которые мы будем ссылаться при описании алгоритма в Разделе 5.

При нахождении планов для каждого отношения, как исходного, так и того, которое встречается в качестве промежуточного результата при выполнении запроса, оптимизатор сохраняет несколько планов. Планы идентифицируются ключами (структура `PathKey`), которые содержат информацию о порядке записей, получаемом при выполнении данного плана. При добавлении очередного плана в список планов для текущего отношения, он сравнивается с другим планом с таким же ключём. Если такой план имеется и его стоимость ниже, то новый план отбрасывается, а иначе новый план заменяет старый.

-
1. Построение планов для всех ОТВ.
 - Вызов `subquery_planner` для каждого ОТВ.
 2. Преобразование `ANY` и `EXISTS` в `JOIN`.
 3. Встраивание функций.
 4. «Поднятие» подзапросов — рассматривается возможность объединения подзапроса с его родителем.
 5. Предобработка, упрощение выражений в разных частях запроса.
 - Вызов `subquery_planner` для подзапросов в `EXISTS` и `ANY`.
 6. Перемещение условий из `HAVING` в `WHERE`.
 7. Вызов `grouping_planner`
 - 7.1. Вызов `query_planner` для нахождения планов выполнения соединений.
 - Вызов `subquery_planner` для подзапросов в `FROM`
 - 7.2. Нахождение планов группировки.
 - 7.3. Нахождение планов для оконных функций.
 - 7.4. Нахождение планов для `DISTINCT`
 - 7.5. Нахождение планов сортировки. (Обработка `ORDER BY`)
 - 7.6. Добавление к планам узла для `LIMIT`
 8. Выбор самого дешёвого плана.
-

Таблица 2: Упрощённая последовательность действий при планировании запроса. Содержание вызова `query_planner`.

5. Возможность реализации алгоритма встраивания ОТВ в PostgreSQL

Рассмотрим, как можно реализовать способ оптимизации, основанный на встраивании, который был описан в Разделе 2, в оптимизаторе PostgreSQL.

Перед тем как встраивать ОТВ, нужно убедиться, что это не изменит семантику запроса. В случае PostgreSQL, это означает, что подзапрос

- не является «модифицирующим», то есть имеет вид `SELECT ...`,
- не использует `VOLATILE` или `STABLE` функции.

Если запрос, удовлетворяет этим условиям, то мы можем рассматривать его встраивание.

Начнём с простого случая, когда имеется только одна ссылка на ОТВ. Тогда мы можем встроить запрос без проверки стоимости. Это можно сделать в самом начале функции `subquery_planner` (Таблица 2) перед построением планов для ОТВ (пункт 1). Так как все оптимизации выполняются после, то нам не нужно ничего делать дополнительно, только удалить его из списка ОТВ, которые лежат в структуре `PlannerInfo` для корневого запроса.

Теперь рассмотрим общий случай. Для его реализации нам нужно понять

1. как добавить в рассмотрение альтернативные планы, включающие встраивание ОТВ,
2. как отложить сравнение стоимости планов до наименьшего общего предка?

Разберём представленные вопросы по-порядку. Начнём с краткого описания решения первой проблемы в оптимизаторе `Orca` [9]. В нём добавление альтернативных планов осуществляется путём применения правил трансформации к элементам, содержащимся в структуре `Memo` [6].

Мемо — это ациклический граф из групп. Каждая группа представляет собой логический оператор и содержит все возможные его реализации (физический оператор или их последовательность), которые генерируются правилами трансформации.

В PostgreSQL единицей, содержащей информацию о логически эквивалентных операциях, является структура `RelOptInfo`. Она представляет собой абстракцию некоторого отношения для оптимизатора. Кроме прочей информации, `RelOptInfo` содержит список путей (`Path`) с помощью которого это отношение можно получить. Структура `Path` и её производные, по сути, описывают различные физические операторы. Дерево из структур `Path` мы будем называть путём или планом.

Таким образом, чтобы добавить альтернативные планы, нам нужно в `RelOptInfo`, которое содержит информацию об отношении, получаемом по ссылке на ОТВ, добавить пути `SubqueryScanPath`, которые описывают получение этого отношения, как результат обращения к подзапросу. Однако мы не можем так поступить, так как `RelOptInfo` содержит информацию для оптимизатора, которая будет различна для ссылки на ОТВ и подзапроса.

Возможным решением может быть создание специальной структуры, аналогичной `AlternativeSubPlan`, которая представляла бы альтернативу между несколькими `RelOptInfo`. Это позволило бы применять при оптимизации методы, которые требуют сравнения альтернатив, связанных с возможной трансформацией запроса.

Перейдём ко второму вопросу. Как уже упоминалось в Разделе 4, при выборе путей, оптимизатор сохраняет самый дешёвый путь с учётом получающегося порядка записей, а также отдельно сохраняет самый дешёвый путь без учёта порядка и путь с минимальной стартовой стоимостью (стоимость до получения первого кортежа ответа). Однако, как обсуждалось в Разделе 2, мы не можем сразу учесть полную стоимость для путей, в которых используются встроенные ОТВ. Для этого нам нужно сохранить их до наименьшего общего предка ссылок на ОТВ. По всей видимости, мы должны выделить отдельный список для таких путей и сравнивать их отдельно между собой, если их ключи

равны.

Таким образом, для реализации оптимизации нетривиального случая, потребуется довольно серьёзный рефакторинг оптимизатора и, скорей всего, добавление новых структур данных.

5.1. Сравнение планов со встроенными ОТВ

Обсудим подробнее как можно сравнивать пути со встроенными ОТВ. Мы рассматриваем пути сгенерированные для получения одного и того же отношения, имеющие равные ключи путей (PathKey). Это означает, что с точки зрения планировщика они взаимозаменяемы.

Более формально отношение взаимозаменяемости путей можно описать следующим образом. Пусть существуют пути p и q . Каждый путь представляет собой дерево из операций. Будем обозначать $p \leftarrow q$, если путь q является поддеревом p . Тогда пути p_i и p_j называются взаимозаменяемыми, если для любого пути $p \leftarrow p_i$, существует путь $p \leftarrow p_j$, и наоборот — для любого пути $p \leftarrow p_j$ существует путь $p \leftarrow p_i$.

Пусть мы хотим сравнить два взаимозаменяемых плана, которые имеют операции обращения к материализованному или встроенному ОТВ. Пусть всего в запросе есть M различных ОТВ. При этом нв ОТВ M_i в запросе m_i ссылок. Обозначим также

- n_i^p — количество сканирований материализованного M_i , используемое в плане p ,
- r_i^p — количество сканирований материализованного M_i , используемое в некотором *родительском* плане p ,
- R — количество ссылок на ОТВ в остальной части запроса (часть, которая не покрывается рассматриваемыми путями),
- C_i — цена материализации M_i ,
- C^p — цена плана p без учёта стоимости материализации используемых в нём ОТВ.

Предположим мы хотим сравнить два взаимозаменяемых плана p_1 и p_2 . Стоимость плана p_j с учётом материализации, при некотором родительском пути можно выразить как:

$$\text{cost}(p_j; r_1, \dots, r_M) = C^{p_j} + \sum_{i=1}^M \text{addition_cost}(p_j, r_i) \quad (1)$$

где $\text{addition_cost}(p_j, r_i)$ — доля стоимости материализации M_i , которая распределяется на путь p_j . При этом, r_i задаются некоторым родительским планом p .

$$\text{addition_cost}(p_j, r_i) = \begin{cases} \frac{C_i \cdot n_i^{p_j}}{n_i^{p_j} + r_i}, & \text{если } n_i^{p_j} > 0 \\ 0, & \text{иначе} \end{cases} \quad (2)$$

Тогда мы можем сравнить стоимость планов с помощью неравенства, получающегося из равенств 1 и 2:

$$C^{p_1} - C^{p_2} < \sum_{i=1}^M \frac{C_i \cdot r_i \cdot (n_i^{p_2} - n_i^{p_1})}{(n_i^{p_1} + r_i)(n_i^{p_2} + r_i)} \quad (3)$$

Ясно, что если неравенство будет выполняться при минимальном значении суммы справа, то мы должны оставить только план p_1 . В противном случае, какой путь выгоднее будет зависеть от пути, выбранного для остальной части запроса, поэтому мы должны оставить оба кандидата.

Рассмотрим, при каких r_i правая сумма принимает минимальное значения. Каждое слагаемое суммы мы можем рассматривать независимо. Для удобства обозначим его как:

$$\text{cost_diff}(r_i) = \frac{C_i \cdot r_i \cdot (n_i^{p_2} - n_i^{p_1})}{(n_i^{p_1} + r_i)(n_i^{p_2} + r_i)} \quad (4)$$

Будем по-отдельности рассматривать случаи различных значений $n_i^{p_2}$ и $n_i^{p_1}$.

$n_i^{p_2} = n_i^{p_1}$ Это тривиальный случай, так как тогда $\text{cost_diff}(r_i) \equiv 0$.

$n_i^{p_1} = 0$ В этом случае

$$\text{cost_diff}(r_i) = \frac{C_i \cdot n_i^{p_2}}{n_i^{p_2} + r_i}$$

и минимум достигается при $r_i = R$. Таким образом

$$\min_cost_diff(r_i) = \frac{C_i \cdot n_i^{p_2}}{n_i^{p_2} + R} \quad (5)$$

$n_i^{p_2} = 0$ В этом случае

$$\text{cost_diff}(r_i) = -\frac{C_i \cdot n_i^{p_1}}{n_i^{p_1} + r_i}$$

и минимум достигается при $r_i = 0$. Таким образом

$$\min_cost_diff(r_i) = -C_i \quad (6)$$

Перед тем как рассмотреть оставшиеся случаи отметим, что $\text{cost_diff}(r_i) = 0$, при $r_i = 0$, а также, что она непрерывно дифференцируема по r_i на $[0, \infty)$. Её производная представлена в Равенстве 7.

$$\begin{aligned} \text{cost_diff}'(r) &= \\ &= \frac{C(n^{p_2} - n^{p_1})((n^{p_2} + r)(n^{p_1} + r) - r((n^{p_2} + n^{p_1}) + 2r^2))}{(n^{p_2} + r)^2(n^{p_1} + r)^2} \\ &= \frac{C(n^{p_2} - n^{p_1})(n^{p_1}n^{p_2} + r^2 - 2r^3)}{(n^{p_2} + r)^2(n^{p_1} + r)^2} \end{aligned} \quad (7)$$

$n_i^{p_2} > n_i^{p_1} > 0$ Тогда $\text{cost_diff}(r_i) > 0$ на интервале $(0, \infty)$, причём $\text{cost_diff}(r_i) \xrightarrow[r_i \rightarrow \infty]{} +0$. Учитывая непрерывность производной, мы можем сказать, что у неё должен быть минимум один корень на этом интервале.

Проанализируем как меняется знак производной на $[0, \infty)$. Очевидно, что при фиксированных n^{p_1} и n^{p_2} , знак будет зависеть только от члена $(n^{p_1}n^{p_2} + r^2 - 2r^3)$. Запишем соответствующее неравенство $n^{p_1}n^{p_2} + r^2 -$

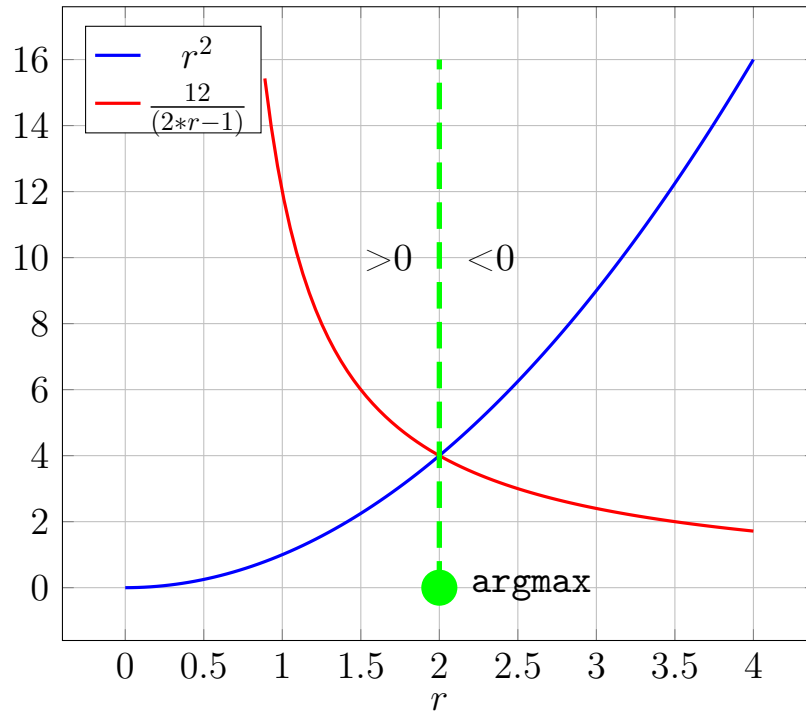


Рис. 2: Иллюстрация смены знака для $r^2 < \frac{n^{p_1}n^{p_2}}{2r-1}$ при $n^{p_1}n^{p_2} = 12$

$2r^3 > 0$ и перепишем его в виде $r^2 < \frac{n^{p_1}n^{p_2}}{2r-1}$. Тогда по графику (Рис. 2) легко понять, что производная на рассматриваемом интервале будет менять знак только один раз. Точка в которой производная равна нулю, в данном случае, будет точкой максимума функции $\text{cost_diff}(r_i)$.

Таким образом, при $n_i^{p_2} > n_i^{p_1} > 0$, минимум будет достигаться в точке $r = 0$, так как в остальных точках $\text{cost_diff}(r_i) > 0$.

$n_i^{p_1} > n_i^{p_2} > 0$ В этом случае $\text{cost_diff}(r_i) \leq 0$ на рассматриваемом отрезке, а значит минимум достигается в точке $\text{argmax cost_diff}'(r)$. Данное значение легко предварительно подсчитать с помощью приведённого выше неравенства для всех разумных (т.е. не слишком больших) значений $n^{p_1}n^{p_2}$.

Сводная информация о значениях на которых достигается минимум

$\text{cost_diff}(r_i)$ представлена в равенстве 8.

$$\text{argmin cost_diff}(r_i) = \begin{cases} 0, & \text{если } n_i^{p_2} > n_i^{p_1} > 0 \\ \text{argmax cost_diff}'(r), & \text{если } n_i^{p_1} > n_i^{p_2} > 0 \\ m_i - R, & \text{если } n_i^{p_1} = 0 \\ 0, & \text{если } n_i^{p_2} = 0 \end{cases} \quad (8)$$

В результате, чтобы сравнить планы, в которых есть обращение к ОТВ, мы можем посчитать минимально возможное значение суммы справа в неравенстве 3. Если при этом неравенство выполняется, то мы оставляем план с меньшей стоимостью. В противном случае мы не можем принять окончательного решения и должны оставить оба плана для дальнейшего рассмотрения.

6. Мнение сообщества разработчиков PostgreSQL о встраивании ОТВ

Предложение о реализации алгоритма встраивания ОТВ было высказано автором в рассылке для разработчиков PostgreSQL [11] в теме «CTE inlining». Проблемы, возникающие из-за того, что оптимизатор считает ОТВ «оптимизационной решёткой», несколько раз поднимались раньше^{2 3 4}, однако широкого обсуждения и принятие какого-либо решения не последовало.

В этот раз в обсуждении участвовало 3 core разработчика из 5, 5 major разработчиков, и ещё несколько контрибьюторов. Участники дискуссии сошлись во мнении, что «оптимизационная решётка» должна быть убрана по-умолчанию, а для возможности использования текущего поведения, необходимо расширить стандартный синтаксис объявления обобщённых выражений опциональным словом `MATERIALIZIZE`, которое можно будет добавить после `WITH` в объявлении ОТВ.

²<https://www.postgresql.org/message-id/flat/4EA6E252.6030002%40linos.es#4EA6E252.6030002@linos.es>

³<https://www.postgresql.org/message-id/flat/29918.1320244719%40sss.pgh.pa.us#29918.1320244719@sss.pgh.pa.us>

⁴<https://www.postgresql.org/message-id/201209191305.44674.db@kavod.com>

Заключение

В данной работе были рассмотрены методы оптимизации, применяемые к обобщённым табличным выражениям в различных СУБД. На основе существующих реализаций был предложен подход к реализации в PostgreSQL встраивания ОТВ, как основного метода оптимизации. Также был описан способ сравнения стоимости путей со встроенными ОТВ при построении плана выполнения запроса. Кроме того, были отмечены недостатки работы оптимизатора PostgreSQL с ОТВ в текущей версии СУБД и приведены примеры запросов, иллюстрирующие их. Предлагаемая инициатива по реализации механизма встраивания ОТВ получила поддержку сообщества разработчиков PostgreSQL.

Список литературы

- [1] Cost-based Query Transformation in Oracle / Rafi Ahmed, Allison Lee, Andrew Witkowski et al. // Proceedings of the 32Nd International Conference on Very Large Data Bases.— VLDB '06.— VLDB Endowment, 2006.— P. 1026–1036.— URL: <http://dl.acm.org/citation.cfm?id=1182635.1164215>.
- [2] Dayal Umeshwar. Of Nests and Trees : A Unified Approach to Processing Queries That Contain Nested Subqueries, Aggregates, and Quantifiers // Vldb.— 1987.— P. 197–208.
- [3] Efficient exploitation of similar subexpressions for query processing / Jingren Zhou, Per Ake Larson, Johann Christoph Freytag, Wolfgang Lehner // SIGMOD '07: Proceedings of the 2007 ACM SIGMOD international conference on Management of data.— 2007.— no. April.— P. 533–544.— URL: <http://portal.acm.org/citation.cfm?id=1247480.1247540>.
- [4] Enhanced subquery optimizations in Oracle / Srikanth Bellamkonda, Rafi Ahmed, Andrew Witkowski et al. // Proceedings of the VLDB Endowment.— 2009.— aug.— Vol. 2, no. 2.— P. 1366–1377.— URL: <http://dl.acm.org/citation.cfm?doid=1687553.1687563>.
- [5] Execution strategies for SQL subqueries / Mostafa Elhemali, César a. Galindo-Legaria, Torsten Grabs, Milind M. Joshi // Proceedings of the 2007 ACM SIGMOD international conference on Management of data - SIGMOD '07.— 2007.— P. 993.— URL: <http://portal.acm.org/citation.cfm?doid=1247480.1247598>.
- [6] Graefe Goetz. The Cascades framework for query optimization // Data Engineering Bulletin.— 1995.— Vol. 18, no. 3.— P. 19–29.— URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.98.9460{&}rep=rep1{&}type=pdf>.
- [7] Kim Won. On optimizing an SQL-like nested query // ACM

Transactions on Database Systems. — 1982. — Vol. 7, no. 3. — P. 443–469.

- [8] Optimization of common table expressions in MPP database systems / Amr El-Helw, Venkatesh Raghavan, Mohamed A. Soliman et al. // Proceedings of the VLDB Endowment. — 2015. — Vol. 8, no. 12. — P. 1704–1715. — URL: <http://dl.acm.org/citation.cfm?id=2824032.2824068>.
- [9] Orca: A modular query optimizer architecture for big data // Sigmod. — 2014. — P. 337–348. — URL: <http://dl.acm.org/citation.cfm?doid=2588555.2595637>.
- [10] PostgreSQL repository // The PostgreSQL Licence (PostgreSQL), The PostgreSQL Global Development Group. — URL: <https://github.com/postgres/postgres> (дата обращения: 02.05.2017).
- [11] postgresql-hackers mailing list, The PostgreSQL Global Development Group. — URL: <https://www.postgresql.org/list/postgresql-hackers/> (дата обращения: 22.05.2017).
- [12] Демонстрационная база данных // The PostgreSQL Licence (PostgreSQL), Postgres Professional. — URL: <https://postgrespro.ru/education/demodb> (дата обращения: 02.05.2017).

Список Листингов

1.	Количество мест в самолёте, с учётом класса обслуживания.	9
1.	План к Запросу 1.	9
2.	Запрос с ОТВ, эквивалентный по результату Запросу 1.	9
2.	План к Запросу 2.	9
3.	Нахождение стыковочных рейсов из Домодедово в Пулково.	11
3.	План выполнения Запроса 3.	12
4.	Нахождение стыковочных рейсов, со встроенными FlightSeatsAvailable и FlightSeatsOccupied.	14
4.	План выполнения Запроса 4. Часть 1.	15
4.	План выполнения Запроса 4. Часть 2.	16