



Санкт-Петербургский государственный политехнический университет

Институт компьютерных наук и технологий

Выполнил студент гр.3530904/90102

Мэн Цзянин

Руководитель

Прокофьев О.В.

Санкт-Петербург

2022

1. Лабораторная работа No.3

1.1 Проектирование аналитической схемы базы данных

1.1.1 Постановка задачи

1.1.2 Реализация

1.1.2.1 ER диаграммы

1.1.2.2 Хранимая процедура (генератор)

1.1.2.3 Анализ плана выполнения запроса

1.2 Использование документно-ориентированных объектов типа Json

1.2.1 Постановка задачи

1.2.2 Реализация

1.2.2.1 IMDB-JSONB

1.2.2.1.1 Результат

1.2.2.2 TOAST

2. Приложение

2.1 Адрес репозитория GitHub

1. Лабораторная работа No.3

1.1 Проектирование аналитической схемы базы данных

<http://www.postgres.cn/docs/14/plpython-funcs.html>

<https://www.cnblogs.com/whitebai/p/12924270.html>

<http://www.postgres.cn/docs/14/performance-tips.html>

<https://juejin.cn/post/6960674004969455623>

1.1.1 Постановка задачи

Седьмое практическое задание связано с проектированием схемы базы данных для аналитики. Будем исходить из того, что приложение, для которого была сделана база данных в задании стала очень популярной и по ней каждый день можно собирать большой объем статистической информации. Результатом данного практического задания являются: **скрипты создания базы данных, хранимая процедура (генератор) для ее заполнения, анализ плана выполнения запроса.**

Требования к БД

- Как минимум одна таблица должна содержать не меньше 10 млн. записей, которые со временем теряют актуальность.
- Другая таблица, связанная с первой, должна содержать не меньше 1 млн. записей.
- В одной из таблиц с количеством записей больше 1 млн. должна быть колонка с текстом, по которой будет необходимо настроить полнотекстовый поиск.

Практическая часть включает:

1. наполнение таблицы, для этого нужно написать хранимую процедуру - генератор на языке `plpython3u`, которая использует словари (для строковых типов), случайные значения (для строковых, числовых типов).
2. оценку скорости выполнения запросов.

Для этого могут быть использованы механизмы секционирования, наследования и индексов. Необходимо подготовить два запроса:

- Запрос к одной таблице, содержащий фильтрацию по нескольким полям.
- Запрос к нескольким связанным таблицам, содержащий фильтрацию по нескольким полям.

Для каждого из этих запросов необходимо провести следующие шаги:

- Получить план выполнения запроса без использования индексов (удаление индекса или отключение его использования в плане запроса).
- Получить статистику (IO и Time) выполнения запроса без использования индексов.
- Создать нужные индексы, позволяющие ускорить запрос.
- Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

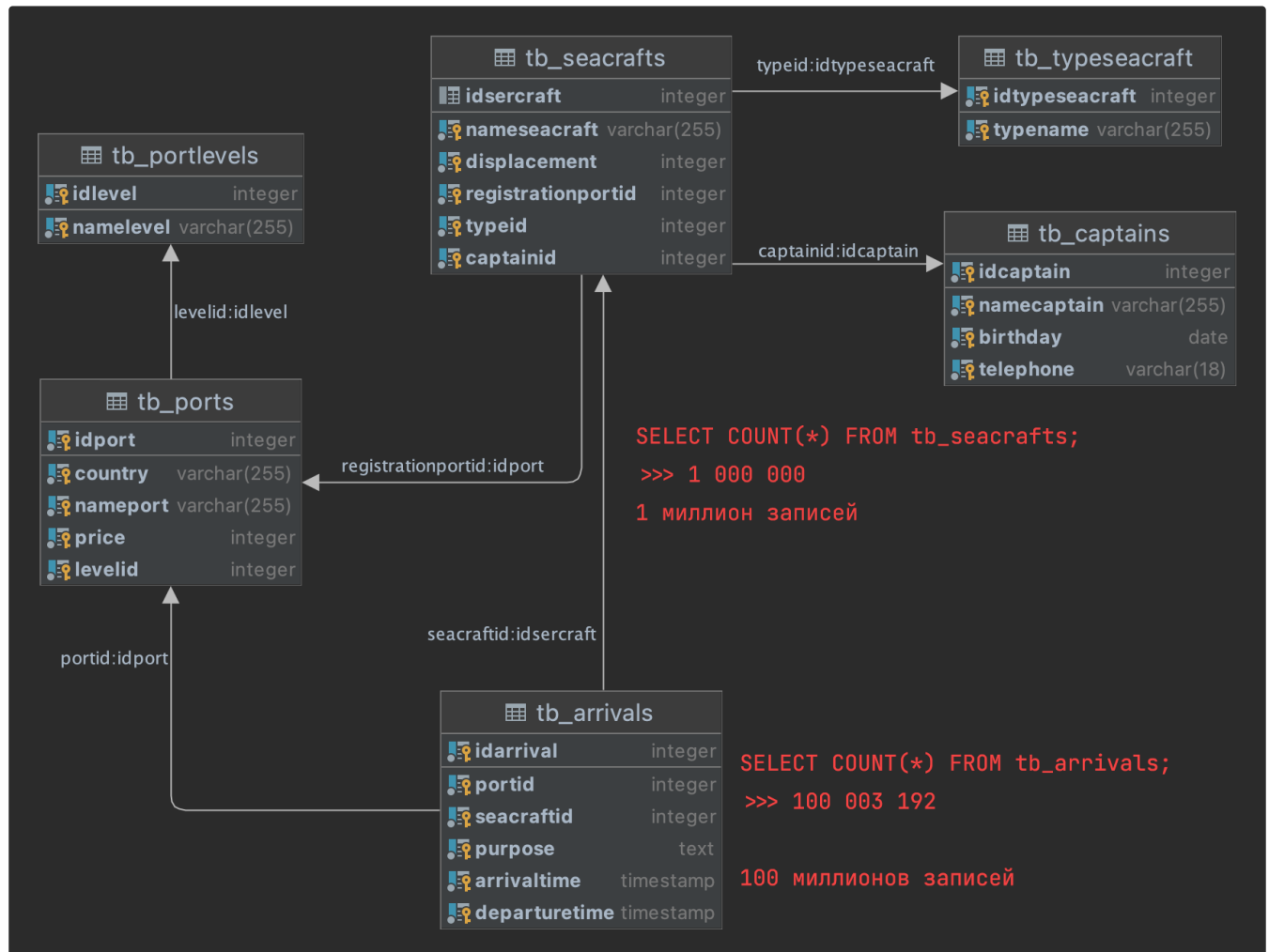
- Получить статистику выполнения запроса с использованием индексов и сравнить с первоначальной статистикой.
- Оценить эффективность выполнения оптимизированного запроса.

Также необходимо продемонстрировать полезность индексов для организации полнотекстового поиска.

Для таблицы объемом 10 млн. записей произвести оптимизацию, позволяющую быстро удалять старые данные, ускорить вставку и чтение данных.

1.1.2 Реализация

1.1.2.1 ER диаграммы



1.1.2.2 Хранимая процедура (генератор)

Скрипты для генерации и вставки случайных записей:

Вставьте 100 000 000 случайных записей данных в tb_arrivals

```

1 CREATE EXTENSION plpython3u;
2
3 -----
4 -- Вставьте 100 000 000 случайных записей данных в tb_arrivals
5 CREATE OR REPLACE FUNCTION fc_InsertArrivals()
6 RETURNS TEXT
7 AS $$
8 import random
9 import datetime
10
11
12 COL_NUM = 100000000
13
14
15 def get_random_text():
16     """
17     Сгенерируйте случайный `текст`.
18
19     :return: text -> string
20     """
21     text = ""
22     num_word = random.randint(1, 20)
23     for i in range(num_word):
24         len_word = random.randint(1, 10)
25         word = "".join(random.sample('zyxwvutsrqponmlkjihgfedcba', len_word))
26         text += (word + ' ')
27
28     return "" + text + ""
29
30
31 def get_random_start_end_time():
32     """
33     Получить строку со случайным начальным временем `start_time_str` и случайным конечным
34     временем `end_time_str`,
35     объединенными вместе
36
37     :return: string -> start time and end time
38     """
39     start = '1900-01-01 00:00:00'
40     end = '2022-04-15 00:00:00'
41     frmt = '%Y-%m-%d %H:%M:%S'
42     num_hour = random.randint(12, 7 * 24)
43     stime = datetime.datetime.strptime(start, frmt)
44     etime = datetime.datetime.strptime(end, frmt)

```

```

44     time_datetime = random.random() * (etime - stime) + stime
45     start_time_str = time_datetime.strftime(frmt)
46     end_time_str = (time_datetime + datetime.timedelta(hours=num_hour)).strftime(frmt)
47
48     return "" + start_time_str + ", " + end_time_str + ""
49
50
51 if __name__ == '__main__':
52     try:
53         num_port = plpy.execute("SELECT COUNT(*) FROM tb_ports")[0]['count']
54         num_seacraft = plpy.execute("SELECT COUNT(*) FROM tb_seacrafts")[0]['count']
55         # return str(num_port)
56
57         for i in range(COL_NUM):
58             portid = random.randint(1, num_port)
59             seacraftid = random.randint(1, num_seacraft)
60             purpose = get_random_text()
61             arrival_leave_time = get_random_start_end_time()
62             command = 'INSERT INTO tb_Arrivals (PortID, SeacraftID, Purpose, ArrivalTime, LeaveTime)
VALUES ('
63             command = command + str(portid) + ', ' + str(seacraftid) + ', ' + purpose + ', ' + arrival_leave_time +
';'
64             # return command
65             plpy.execute(command)
66
67     except plpy.SPIError as e:
68         return "[ERROR] %s" % e.sqlstate
69
70     else:
71         return "[Info] Successful insert"
72
73 $$ LANGUAGE plpython3u;
74
75 -----
76
77 TRUNCATE TABLE tb_arrivals CASCADE;
78 SELECT fc_InsertArrivals(); -- 运行时间 13040 s
79
80
81 SELECT COUNT(*) FROM tb_arrivals; -- 运行时间 26.195 s
82

```


Вставьте 1 000 000 случайных записей данных в tb_arrivals

```

1
2 -----
3 -- 向 tb_seacraft 中插入 1000000 条随机数据记录
4 -- Вставьте 1 000 000 случайных записей данных в tb_arrivals
5 CREATE OR REPLACE FUNCTION fc_InsertSeacrafts()
6     RETURNS TEXT
7     AS $$
8 import random
9 import datetime
10
11
12 COL_NUM = 1000000
13
14
15 def get_random_name():
16     """
17     生成一段随机 `name`
18     Сгенерируйте случайный `name`.
19
20     :return: text -> string
21     """
22     len_word = random.randint(1, 10)
23     word = ".join(random.sample('zyxwvutsrqponmlkjihgfedcba', len_word))
24
25     return word
26
27
28 if __name__ == '__main__':
29     try:
30         num_captain = plpy.execute("SELECT COUNT(*) FROM tb_captains")[0]['count']
31         num_port = plpy.execute("SELECT COUNT(*) FROM tb_ports")[0]['count']
32         num_type = plpy.execute("SELECT COUNT(*) FROM tb_portlevels")[0]['count']
33
34         for i in range(COL_NUM):
35             regportid = random.randint(1, num_port)
36             captainid = random.randint(1, num_captain)
37             typeid = random.randint(1, num_type)
38             displacement = random.randint(10000, 1000000)
39             namesecraft = get_random_word()
40
41             command = 'INSERT INTO tb_seacraft(NameSeacraft, Displacement, RegPortID, TypeID,
42             CaptainID) VALUES ('
43             command = command + str(namesecraft) + ', ' + str(displacement) + ', ' + str(regportid) + ', ' +
44             str(typeid) + str(captainid) + ');'
45             # return command
46             plpy.execute(command)
47
48         except plpy.SPIError as e:
49             return "[ERROR] %s" % e.sqlstate

```

```
49     else:
50         return "[Info] Successful insert"
51
52 $$ LANGUAGE plpython3u;
53
54 -----
```

1.1.2.3 Анализ плана выполнения запроса

<http://www.postgres.cn/docs/14/performance-tips.html>

<http://www.postgres.cn/docs/14/using-explain.html#USING-EXPLAIN-BASICS>

<http://www.postgres.cn/docs/14/runtime-config-query.html#RUNTIME-CONFIG-QUERY-CONSTANTS>

<https://blog.csdn.net/kmblack1/article/details/80761647>

Константы стоимости для планировщика

Константы	Описание	Значение по умолчанию (стоимость)
<code>cpu_tuple_cost(floating point)</code>	Задаёт приблизительную стоимость обработки <u>каждой строки</u> при выполнении запроса	0.01
<code>cpu_index_tuple_cost(floating point)</code>	Задаёт приблизительную стоимость обработки <u>каждой записи индекса</u> при сканировании индекса	0.005
<code>cpu_operator_cost(floating point)</code>	Задаёт приблизительную стоимость обработки <u>оператора или функции</u> при выполнении запроса	0.0025

[1] Запрос к одной таблице, содержащий фильтрацию по нескольким полям.

1. Получить план выполнения запроса **без использования индексов** (удаление индекса или отключение его использования в плане запроса).

```

1 | SELECT * FROM tb_arrivals WHERE ArrivalTime BETWEEN '2000-01-01 00:00:00'::TIMESTAMP
2 |   AND '2001-01-01 00:00:00' AND PortID > 1000;
3 | DROP INDEX IX_ArriverTime;
4 | DROP INDEX IX_PortID;
```

2. Получить статистику (IO и Time) выполнения запроса **без использования индексов**.

```

1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE ArrivalTime BETWEEN
  '2000-01-01 00:00:00'::TIMESTAMP AND '2001-01-01 00:00:00' AND PortID > 1000;
2
3
4
5
6
7
8
9
10
11
12

```

QUERY PLAN

```

4 Gather (cost=1000.00..2348941.17 rows=326215 width=97) (actual time=1.252..27359.813
rows=332165 loops=1)
5   Workers Planned: 2
6   Workers Launched: 2
7   -> Parallel Seq Scan on tb_arrivals (cost=0.00..2315319.67 rows=135923 width=97) (actual
time=0.427..27321.096 rows=110722 loops=3)
8     Filter: ((arrivaltime >= '2000-01-01 00:00:00'::timestamp without time zone) AND (arrivaltime <=
'2001-01-01 00:00:00'::timestamp without time zone) AND (portid > 1000))
9     Rows Removed by Filter: 33222612
10    Planning Time: 3.485 ms
11    Execution Time: 27368.246 ms
12    (8 rows)

```

3. Создать нужные индексы, позволяющие ускорить запрос.

```

1 CREATE INDEX IX_ArriverTime ON tb_arrivals(ArrivalTime);
2 CREATE INDEX IX_PortID ON tb_arrivals(PortID);

```

Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```

1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE ArrivalTime BETWEEN
  '2000-01-01 00:00:00'::TIMESTAMP AND '2001-01-01 00:00:00' AND PortID > 1000;
2
3
4
5
6
7
8
9
10
11
12
13
14
15

```

QUERY PLAN

```

4 Gather (cost=459271.16..1940524.67 rows=326215 width=97) (actual time=557.093..23659.834
rows=332165 loops=1)
5   Workers Planned: 2
6   Workers Launched: 2
7   -> Parallel Bitmap Heap Scan on tb_arrivals (cost=458271.16..1906903.17 rows=135923 width=97)
(actual time=537.187..23605.583 rows=110722 loops=3)
8     Recheck Cond: ((arrivaltime >= '2000-01-01 00:00:00'::timestamp without time zone) AND
(arrivaltime <= '2001-01-01 00:00:00'::timestamp without time zone) AND (portid > 1000))
9     Rows Removed by Index Recheck: 12492294
10     Heap Blocks: exact=14079 lossy=199261
11     -> BitmapAnd (cost=458271.16..458271.16 rows=326215 width=0) (actual
time=551.432..551.432 rows=0 loops=1)
12       -> Bitmap Index Scan on ix_arrivertime (cost=0.00..16861.76 rows=804919 width=0) (actual
time=96.028..96.028 rows=818208 loops=1)
13         Index Cond: ((arrivaltime >= '2000-01-01 00:00:00'::timestamp without time zone) AND
(arrivaltime <= '2001-01-01 00:00:00'::timestamp without time zone))
14       -> Bitmap Index Scan on ix_portid (cost=0.00..441246.04 rows=40527663 width=0) (actual
time=448.839..448.839 rows=40516196 loops=1)
15         Index Cond: (portid > 1000)

```

16 Planning Time: 0.142 ms
 17 Execution Time: 23671.716 ms
 18 (14 rows)

4. Получить статистику выполнения запроса с использованием индексов и сравнить с первоначальной статистикой.

	cost	actual time	Planning Time	Execution Time
без использования индексов	1000.00..2348941.17	1.252..27359.813	3.485 ms	27368.246 ms
с использованием индексов	459271.16..1940524.67	557.093..23659.834	0.142 ms	23671.716 ms
Разница	-866687.6599	-4255.82	-3.343 ms	-3696.5299 ms

5. Оценить эффективность выполнения оптимизированного запроса.

После использования индекса (приблизительная стоимость запуска(cost)) значительно увеличилась, однако (приблизительная общая стоимость (cost)) и (фактическое время (Execution Time)) значительно уменьшились. (Фактическое время (Execution Time)) сократилось на 13,506%

[2] Запрос к нескольким связанным таблицам, содержащий фильтрацию по нескольким полям

1. Получить план выполнения запроса **без использования индексов** (удаление индекса или отключение его использования в плане запроса).

```

1  SELECT * FROM tb_arrivals INNER JOIN tb_seacrafts ON
   tb_arrivals.seacraftID=tb_seacrafts.IDSeacraft INNER JOIN tb_ports ON
   tb_arrivals.portID=tb_ports.IDPort WHERE NameSeacraft='cmari' AND NamePort='St Petersburg';
2
3  DROP INDEX IX_NameSeacraft;
4  DROP INDEX IX_NamePort;

```

2. Получить статистику (IO и Time) выполнения запроса **без использования индексов**.

```

1  db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals INNER JOIN tb_seacrafts ON
   tb_arrivals.seacraftID=tb_seacrafts.IDSeacraft INNER JOIN tb_ports ON
   tb_arrivals.portID=tb_ports.IDPort WHERE NameSeacraft='cmari' AND NamePort='St Petersburg';
2
3  QUERY PLAN
4  -----
5  Nested Loop (cost=1008.31..2140204.13 rows=1 width=161) (actual time=698.274..29483.075 rows=6
6  loops=1)
7    Join Filter: (tb_arrivals.seacraftid = tb_seacrafts.idseacraft)
8    Rows Removed by Join Filter: 59539
9    -> Seq Scan on tb_seacrafts (cost=0.00..20033.00 rows=1 width=29) (actual time=0.177..55.336
10   rows=1 loops=1)
11      Filter: ((nameseacraft)::text = 'cmari'::text)
12      Rows Removed by Filter: 999999
13    -> Gather (cost=1008.31..2119427.53 rows=59488 width=132) (actual time=4.875..29424.577
14   rows=59545 loops=1)
15      Workers Planned: 2
16      Workers Launched: 2
17      -> Hash Join (cost=8.31..2112478.73 rows=24787 width=132) (actual time=6.244..29400.751
18   rows=19848 loops=3)
19         Hash Cond: (tb_arrivals.portid = tb_ports.idport)
20         -> Parallel Seq Scan on tb_arrivals (cost=0.00..2002819.67 rows=4166667 width=97) (actual
21   time=0.345..27659.885 rows=3333333 loops=3)
22         -> Hash (cost=8.29..8.29 rows=1 width=35) (actual time=0.165..0.165 rows=1 loops=3)
23             Buckets: 1024 Batches: 1 Memory Usage: 9kB
24         -> Index Scan using uq_ports_nameport on tb_ports (cost=0.28..8.29 rows=1 width=35)
25   (actual time=0.118..0.119 rows=1 loops=3)
26             Index Cond: ((nameport)::text = 'St Petersburg'::text)
27   Planning Time: 7.004 ms
28   Execution Time: 29483.181 ms
29   (18 rows)

```

3. Создать нужные индексы, позволяющие ускорить запрос.

```

1 CREATE INDEX IX_NameSeacraft ON tb_seacrafts(NameSeacraft);
2 CREATE INDEX IX_NamePort ON tb_ports(NamePort);

```

Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```

1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals INNER JOIN tb_seacrafts ON
  tb_arrivals.seacraftID=tb_seacrafts.IDSeacraft INNER JOIN tb_ports ON
  tb_arrivals.portID=tb_ports.IDPort WHERE NameSeacraft='cmari' AND NamePort='St Petersburg';
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22

```

QUERY PLAN

```

-----
 Gather (cost=1016.76..2113212.11 rows=1 width=161) (actual time=9437.606..29046.762 rows=6
 loops=1)
   Workers Planned: 2
   Workers Launched: 2
   -> Hash Join (cost=16.76..2112212.01 rows=1 width=161) (actual time=6601.278..29032.137 rows=2
 loops=3)
     Hash Cond: (tb_arrivals.portid = tb_ports.idport)
     -> Hash Join (cost=8.46..2112203.59 rows=42 width=126) (actual time=6.039..29030.866
 rows=3333 loops=3)
       Hash Cond: (tb_arrivals.seacraftid = tb_seacrafts.idseacraft)
       -> Parallel Seq Scan on tb_arrivals (cost=0.00..2002819.67 rows=4166667 width=97) (actual
 time=0.461..27245.982 rows=3333333 loops=3)
         -> Hash (cost=8.44..8.44 rows=1 width=29) (actual time=0.353..0.356 rows=1 loops=3)
           Buckets: 1024 Batches: 1 Memory Usage: 9kB
           -> Index Scan using ix_nameseacraft on tb_seacrafts (cost=0.42..8.44 rows=1 width=29)
 (actual time=0.349..0.350 rows=1 loops=3)
             Index Cond: ((nameseacraft)::text = 'cmari'::text)
       -> Hash (cost=8.29..8.29 rows=1 width=35) (actual time=0.222..0.223 rows=1 loops=3)
         Buckets: 1024 Batches: 1 Memory Usage: 9kB
         -> Index Scan using ix_nameport on tb_ports (cost=0.28..8.29 rows=1 width=35) (actual
 time=0.215..0.215 rows=1 loops=3)
           Index Cond: ((nameport)::text = 'St Petersburg'::text)
     Planning Time: 3.746 ms
     Execution Time: 29047.011 ms
 (18 rows)

```

4. Создать нужные индексы, позволяющие ускорить запрос.

```

1 CREATE EXTENSION pg_trgm;
2 CREATE INDEX IX_NameSeacraft ON tb_seacrafts USING GIN(NameSeacraft gin_trgm_ops);
3 CREATE INDEX IX_NamePort ON tb_ports USING GIN(NamePort gin_trgm_ops);

```

Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```

1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals INNER JOIN tb_seacrafts ON
  tb_arrivals.seacraftID=tb_seacrafts.IDSeacraft INNER JOIN tb_ports ON
  tb_arrivals.portID=tb_ports.IDPort WHERE NameSeacraft='cmari' AND NamePort='St Petersburg';

```

```

2                                     QUERY PLAN
3 -----
4 Gather (cost=1088.34..2321635.00 rows=1 width=161) (actual time=11391.572..35022.070 rows=6
5 loops=1)
6   Workers Planned: 2
7   Workers Launched: 2
8   -> Hash Join (cost=88.34..2320634.90 rows=1 width=161) (actual time=12996.658..35002.537
9 rows=2 loops=3)
10     Hash Cond: (tb_arrivals.portid = tb_ports.idport)
11     -> Hash Join (cost=80.03..2320626.48 rows=42 width=126) (actual time=10.302..34999.074
12 rows=3333 loops=3)
13       Hash Cond: (tb_arrivals.seacraftid = tb_seacrafts.idseacraft)
14       -> Parallel Append (cost=0.00..2211170.01 rows=41667034 width=97) (actual
15 time=1.530..32921.062 rows=33333333 loops=3)
16       -> Parallel Seq Scan on tb_arrivals tb_arrivals_1 (cost=0.00..2002819.67 rows=41666667
17 width=97) (actual time=1.517..31469.568 rows=33333333 loops=3)
18       -> Parallel Seq Scan on tb_child_arrivals tb_arrivals_2 (cost=0.00..15.18 rows=518
19 width=60) (actual time=0.001..0.001 rows=0 loops=1)
20       -> Hash (cost=80.02..80.02 rows=1 width=29) (actual time=1.348..1.349 rows=1 loops=3)
21         Buckets: 1024 Batches: 1 Memory Usage: 9kB
22         -> Bitmap Heap Scan on tb_seacrafts (cost=76.01..80.02 rows=1 width=29) (actual
23 time=1.346..1.347 rows=1 loops=3)
24           Recheck Cond: ((name_seacraft)::text = 'cmari'::text)
25           Heap Blocks: exact=1
26           -> Bitmap Index Scan on ix_name_seacraft (cost=0.00..76.01 rows=1 width=0) (actual
27 time=1.339..1.339 rows=1 loops=3)
28             Index Cond: ((name_seacraft)::text = 'cmari'::text)
29       -> Hash (cost=8.29..8.29 rows=1 width=35) (actual time=0.151..0.152 rows=1 loops=3)
30         Buckets: 1024 Batches: 1 Memory Usage: 9kB
31         -> Index Scan using uq_ports_nameport on tb_ports (cost=0.28..8.29 rows=1 width=35)
32 (actual time=0.083..0.084 rows=1 loops=3)
33       Index Cond: ((nameport)::text = 'St Petersburg'::text)
34 Planning Time: 1.025 ms
35 Execution Time: 35022.279 ms
36 (23 rows)

```

5. Получить статистику выполнения запроса с использованием индексов и сравнить с первоначальной статистикой.

	cost	actual time	Planning Time	Execution Time	Разница Execution Time с неиспользуемыми индексами
без использования индексов	1008.31..2140204.13	698.274..29483.075	7.004 ms	29483.181 ms	
с использованием индексов	1016.76..2113212.11	9437.606..29046.762	3.746 ms	29047.011 ms	-436.17 ms
с использованием индексов (GIN)	1088.34..2321635.00	11180.878..35380.249	1.025	35022.279 ms	+5539.09 ms

6. Оценить эффективность выполнения оптимизированного запроса.

При добавлении индексов к двум полям **varchar** подтаблицы значительного улучшения **в скорости запросов не наблюдается**. GIN-индексирование в этом случае вместо того, чтобы ускорить запрос, делает его медленнее.

[3] Также необходимо продемонстрировать полезность индексов для организации полнотекстового поиска.

1. Получить план выполнения запроса **без использования индексов** (удаление индекса или отключение его использования в плане запроса).

```
1 | SELECT * FROM tb_arrivals WHERE Purpose LIKE '%at%';
2 |
3 | DROP INDEX IX_Purpose;
```

2. Получить статистику (IO и Time) выполнения запроса **без использования индексов**.

```
1 | db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE Purpose LIKE '% at %';
2 |                                     QUERY PLAN
3 | -----
4 | Gather (cost=1000.00..2108981.63 rows=9953 width=97) (actual time=2.708..35207.244 rows=145614
   | loops=1)
5 |   Workers Planned: 2
6 |   Workers Launched: 2
7 |   -> Parallel Seq Scan on tb_arrivals (cost=0.00..2106986.33 rows=4147 width=97) (actual
   | time=2.798..35181.150 rows=48538 loops=3)
8 |     Filter: (purpose ~ '% at %':text)
9 |     Rows Removed by Filter: 33284795
10 | Planning Time: 0.337 ms
11 | Execution Time: 35212.464 ms
12 | (8 rows)
```

3. Создать нужные индексы, позволяющие ускорить запрос.

```
1 | CREATE INDEX IX_Purpose ON tb_arrivals(Purpose);
```

Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```

1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE Purpose LIKE '% at %';
2                                     QUERY PLAN
3 -----
4 Gather (cost=1000.00..2108981.63 rows=9953 width=97) (actual time=3.036..36736.283 rows=145614
loops=1)
5   Workers Planned: 2
6   Workers Launched: 2
7   -> Parallel Seq Scan on tb_arrivals (cost=0.00..2106986.33 rows=4147 width=97) (actual
time=0.982..36713.617 rows=48538 loops=3)
8     Filter: (purpose ~ '% at %'::text)
9     Rows Removed by Filter: 33284795
10  Planning Time: 5.599 ms
11  Execution Time: 36741.591 ms
12  (8 rows)

```

4. Создать нужные индексы (GIN), позволяющие ускорить запрос.

```

1 CREATE INDEX IX_Purpose ON tb_arrivals USING GIN(Purpose gin_trgm_ops);

```

Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```

1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE Purpose LIKE '% at %';
2                                     QUERY PLAN
3 -----
4 Append (cost=1277.14..38806.13 rows=9960 width=97) (actual time=711.662..21840.648
rows=145614 loops=1)
5   -> Bitmap Heap Scan on tb_arrivals tb_arrivals_1 (cost=1277.14..38735.33 rows=9953 width=97)
(actual time=711.661..21829.443 rows=145614 loops=1)
6     Recheck Cond: (purpose ~ '% at %'::text)
7     Rows Removed by Index Recheck: 8520120
8     Heap Blocks: exact=33655 lossy=137020
9     -> Bitmap Index Scan on ix_purpose (cost=0.00..1274.65 rows=9953 width=0) (actual
time=706.322..706.322 rows=180530 loops=1)
10       Index Cond: (purpose ~ '% at %'::text)
11   -> Seq Scan on tb_child_arrivals tb_arrivals_2 (cost=0.00..21.00 rows=7 width=60) (actual
time=0.062..0.062 rows=0 loops=1)
12     Filter: (purpose ~ '% at %'::text)
13  Planning Time: 27.951 ms
14  Execution Time: 21846.879 ms
15  (11 rows)

```

5. Получить статистику выполнения запроса с использованием индексов и сравнить с первоначальной статистикой.

	cost	actual time	Planning Time	Execution Time	Разница Execution Time с неиспользуемыми индексами
без использования индексов	1000.00..2108981.63	2.708..35207.244	0.337 ms	35212.464 ms	
с использованием индексов	1000.00..2108981.63	3.032..36313.774	4.645 ms	36319.114 ms	+1106.65 ms
с использованием индексов(GIN)	1277.14..38806.13	711.662..21840.648	27.951 ms	21846.879 ms	-13365.585 ms

6. Оценить эффективность выполнения оптимизированного запроса.

Несмотря на то, что для текстового поля создан индекс, база данных **не использует его для запросов**.

При использовании индекса GIN скорость выполнения запросов значительно выше.

[4] Также необходимо продемонстрировать полезность индексов для организации полнотекстового поиска.

1. Получить план выполнения запроса **без использования индексов** (удаление индекса или отключение его использования в плане запроса).

```
1 SELECT * FROM tb_arrivals WHERE char_length(Purpose) > 20;
2
3 DROP INDEX IX_Purpose;
```

2. Получить статистику (IO и Time) выполнения запроса **без использования индексов**.

```
1 db_port=## EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE char_length(Purpose) > 20;
2 QUERY PLAN
3 -----
4 Append (cost=0.00..3252844.33 rows=33333626 width=97) (actual time=0.522..22543.607
rows=86235207 loops=1)
5   -> Seq Scan on tb_arrivals tb_arrivals_1 (cost=0.00..3086153.00 rows=33333333 width=97) (actual
time=0.521..19450.968 rows=86235207 loops=1)
6     Filter: (char_length(purpose) > 20)
7     Rows Removed by Filter: 13764793
8   -> Seq Scan on tb_child_arrivals tb_arrivals_2 (cost=0.00..23.20 rows=293 width=60) (actual
time=0.026..0.026 rows=0 loops=1)
9     Filter: (char_length(purpose) > 20)
10  Planning Time: 2.516 ms
11  Execution Time: 24233.586 ms
12  (8 rows)
```

3. Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```
1 CREATE INDEX IX_Purpose ON tb_arrivals USING GIN(Purpose gin_trgm_ops);
```

Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```

1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE char_length(Purpose) > 20;
2                                     QUERY PLAN
3 -----
4 Append (cost=0.00..3252844.33 rows=33333626 width=97) (actual time=0.514..23314.352
rows=86235207 loops=1)
5   -> Seq Scan on tb_arrivals tb_arrivals_1 (cost=0.00..3086153.00 rows=33333333 width=97) (actual
time=0.513..20128.237 rows=86235207 loops=1)
6     Filter: (char_length(purpose) > 20)
7     Rows Removed by Filter: 13764793
8   -> Seq Scan on tb_child_arrivals tb_arrivals_2 (cost=0.00..23.20 rows=293 width=60) (actual
time=0.042..0.042 rows=0 loops=1)
9     Filter: (char_length(purpose) > 20)
10  Planning Time: 0.602 ms
11  Execution Time: 25058.894 ms
12  (8 rows)

```

4. Получить статистику выполнения запроса с использованием индексов и сравнить с первоначальной статистикой.

	cost	actual time	Planning Time	Execution Time	Разница Execution Time с неиспользуемыми индексами
без использования индексов	0.00..3252844.33	0.522..22543.607	2.516 ms	24233.586 ms	
с использованием индексов(GIN)	0.00..3252844.33	0.514..23314.352	0.602 ms	25058.894 ms	+825.308 ms

5. Оценить эффективность выполнения оптимизированного запроса.

Из полученных результатов видно, что в данном случае индекс GIN не действует (не используется)

[5] Также необходимо продемонстрировать полезность индексов для организации полнотекстового поиска.

1. Получить план выполнения запроса **без использования индексов** (удаление индекса или отключение его использования в плане запроса).

```
1 SELECT * FROM tb_arrivals WHERE Purpose='sfrht';
2
3 DROP INDEX IX_Purpose;
```

2. Получить статистику (IO и Time) выполнения запроса **без использования индексов**.

```
1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE Purpose='sfrht';
2                                     QUERY PLAN
3 -----
4  Gather (cost=1000.00..2108005.45 rows=26 width=91) (actual time=30218.895..30220.101 rows=0
5  loops=1)
6    Workers Planned: 2
7    Workers Launched: 2
8    -> Parallel Append (cost=0.00..2107002.85 rows=10 width=91) (actual time=30192.606..30192.609
9  rows=0 loops=3)
10     -> Parallel Seq Scan on tb_arrivals tb_arrivals_1 (cost=0.00..2106986.33 rows=9 width=97)
11     (actual time=30192.604..30192.604 rows=0 loops=3)
12         Filter: (purpose = 'sfrht'::text)
13         Rows Removed by Filter: 33333333
14     -> Parallel Seq Scan on tb_child_arrivals tb_arrivals_2 (cost=0.00..16.47 rows=2 width=60)
15     (actual time=0.002..0.002 rows=0 loops=1)
16         Filter: (purpose = 'sfrht'::text)
17  Planning Time: 1.095 ms
18  Execution Time: 30220.227 ms
19  (11 rows)
```

3. Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```
1 CREATE INDEX IX_Purpose ON tb_arrivals USING GIN(Purpose gin_trgm_ops);
```

Получить план выполнения запроса с использованием индексов и сравнить с первоначальным планом.

```
1 db_port=# EXPLAIN ANALYZE SELECT * FROM tb_arrivals WHERE Purpose='sfrht';
2                                     QUERY PLAN
3 -----
4  Append (cost=2396.17..2505.33 rows=26 width=91) (actual time=657.419..657.419 rows=0 loops=1)
5    -> Bitmap Heap Scan on tb_arrivals tb_arrivals_1 (cost=2396.17..2484.20 rows=22 width=97) (actual
6    time=657.409..657.409 rows=0 loops=1)
```

```

6      Recheck Cond: (purpose = 'sfrht'::text)
7      Rows Removed by Index Recheck: 23
8      Heap Blocks: exact=23
9      -> Bitmap Index Scan on ix_purpose (cost=0.00..2396.17 rows=22 width=0) (actual
time=642.080..642.080 rows=23 loops=1)
10         Index Cond: (purpose = 'sfrht'::text)
11      -> Seq Scan on tb_child_arrivals tb_arrivals_2 (cost=0.00..21.00 rows=4 width=60) (actual
time=0.008..0.008 rows=0 loops=1)
12         Filter: (purpose = 'sfrht'::text)
13      Planning Time: 14.311 ms
14      Execution Time: 657.908 ms
15      (11 rows)

```

4. Получить статистику выполнения запроса с использованием индексов и сравнить с первоначальной статистикой.

	cost	actual time	Planning Time	Execution Time	Разница Execution Time с неиспользуемыми индексами
без использования индексов	1000.00..2108005.45	30218.895..30220.101	1.095 ms	30220.227 ms	
с использованием индексов(GIN)	2396.17..2505.33	657.419..657.419	14.311 ms	657.908 ms	-29562.319 ms

5. Оценить эффективность выполнения оптимизированного запроса.

В этом случае индекс GIN очень полезен, и скорость выполнения запросов значительно ускоряется

1.2 Использование документно-ориентированных объектов типа Json

1.2.1 Постановка задачи

PostgreSQL стала первой реляционной базой данных, поддерживающей слабоструктурированные данные. В PostgreSQL для этого используется JSON (JavaScript Object Notation, Запись объекта JavaScript RFC 7159), который имеет два представления: json и jsonb. Для реализации эффективного механизма запросов к этим типам данных в Postgres также имеется тип jsonpath. Официально JSON появился в PostgreSQL в 2014 году. PostgreSQL с JSONB совмещает гибкость NoSQL, а также надёжность и богатство функциональности реляционных СУБД.

В практической части необходимо:

- создать БД IMDB test, использующую стандартные атрибуты и атрибут jsonb. Ссылка на интерфейсы. Таблицы находятся на я.диске в папке DataSet. Описание атрибута jsonb:

```
{
  "nconst": "nm0000151",
  "primaryName": "Morgan Freeman",
  "roles": [
    {
      "title": "The Shawshank Redemption",
      "year": "1994",
      "character name": "Ellis Boyd 'Red' Redding"
    },
    {
      "title": "Unforgiven",
      "year": "1992",
      "character name": "Ned Logan"
    },
    {
      "title": "Through the Wormhole",
      "series name": "Are Aliens Inside Us? (#6.5)",
      "year": "2010",
      "character name": "Himself - Narrator"
    }
  ],
  "birthYear": "1937",
  "deathYear": "\N",
}
```

где nconst, birthYear, deathYear это записи таблицы *name.basics.tsv.gz*, roles загружать из таблицы

- Составить 3-4 запроса с использованием jsonb.
- Измерить время доступа к ключу для каждой строчки (в виде таблицы или графика). Оценить влияние длины строки на скорость доступа (линейная зависимость). Как можно это влияние уменьшить.
- Составить запрос на изменение PrimaryName у актёра. Сравнить изменение объёма БД для актёра с малым кол-вом ролей и актёра с большим количеством ролей (toasted roles).

1.2.2 Реализация

<http://www.postgres.cn/docs/14/datatype-json.html>

<https://juejin.cn/post/6844903857009623048>

<https://segmentfault.com/a/1190000019344353>

<https://programmers.buzz/posts/start-exploring-database-indices/>

<https://habr.com/ru/company/oleg-bunin/blog/646987/>

1.2.2.1 IMDB-JSONB

Были использованы все три набора данных: `actors.list.txt`, `actresses.list.txt`, `name.basics.tsv`.

Скрипты, используемые для обработки данных:

1. Из-за того, что одновременная обработка всех данных может занимать много памяти (используя вложенность DataFrame от pandas в качестве структуры данных). Поэтому используйте скрипт Python для разделения набора данных на более мелкие части и их последующего объединения

Объединенный результат сохраняется в виде двоичного файла с использованием сериализации Python для уменьшения использования памяти и жесткого диска, а также для облегчения последующей обработки.



В этом скрипте можно объединить все фильмы одного актера в один

```
1 # -----*----- coding: utf-8 -----*-----
2 # @Time : 2022/4/26 16:39
3 # @Author : 冰糖雪狸 (NekoSilverfox)
4 # @Project : IMDB 数据处理
```

```

5 # @File : 切割提取.py
6 # @Software: PyCharm
7 # @Github : https://github.com/NekoSilverFox
8 # -----
9 import pandas as pd
10 import numpy as np
11 import datetime
12 import pickle
13 import re
14
15
16 def normal_actor(source_data: pd.DataFrame, dump_name_title: str, dump_df_res: str) ->
pd.DataFrame:
17     """
18     Нормализовать данные об актерах IMDB в pd.DataFrame, с отдельной строкой для каждого
актера и их индивидуальных кредитов (аналогично перекрестной таблице)
19     :param file_path:
20     :return: pd.DataFrame
21     """
22
23     """
24     Извлеките название, дату выпуска и название серии с помощью регулярных выражений и
используйте их в качестве нового DataFrame
25     Результатом является DataFrame с MutIndex
26     """
27     print('>>' * 50)
28     print('[INFO] Начните извлекать информацию из строки')
29     time_start = datetime.datetime.now()
30     list_name_title = []
31     i = 1
32     name_list = []
33     for col in source_data.values:
34         if col[0] is not np.nan:
35             names = col[0]
36             name_list = names.split(',')
37
38             if len(name_list) == 0:
39                 continue
40
41             for this_name in name_list:
42
43                 title_mix = col[1]
44
45                 # Регулярные выражения используются для извлечения
46                 """title - Название фильма Первая часть строки """
47                 title = re.search(r'^(\{[*]', title_mix)
48                 if title is not None:
49                     title = str(title.group()[:-1])
50
51                 """year - Год выпуска ()"""
52                 year = re.search(r'(?!=\{1})[d]{4}(?!=\{1})', title_mix)
53                 if year is not None:
54                     year = int(year.group())

```

```

55
56     """series name - Название серии:  {}"""
57     series_name = re.search(r'\{(.*)\}', title_mix)
58     if series_name is not None:
59         series_name = str(series_name.group()[1:-1])
60
61     """character name - Имя персонажа"""
62     character_name = re.search(r'\[(.*)\]', title_mix)
63     if character_name is not None:
64         character_name = str(character_name.group()[1:-1])
65     # name_title.append([this_name, title, series_name, year, character_name])
66     # name_title.append([this_name, [title, series_name, year, character_name]])
67     rols = pd.DataFrame([title, series_name, year, character_name],
68                         columns=['title', 'series name', 'year', 'character name'])
69
70     list_name_title.append([this_name, rols]) # this_name 是 str, rols 是DataFrame
71
72     if i % 10000 == 0:
73         use_sec = (datetime.datetime.now() - time_start).seconds
74         print('[INFO] Обработано ', i, ' строк | ', (i / source_data.shape[0]) * 100, '% | Затраченное
75 время: ', use_sec, ' s (', use_sec / 60,
76         ') min')
77         i += 1
78
79     use_sec = (datetime.datetime.now() - time_start).seconds
80     print('[INFO] Окончание извлечения данных, время: ', use_sec, ' s (', use_sec / 60, ') min')
81
82     """Получить результат df_name_title в виде DataFrame с MutIndex"""
83     print('>>' * 50)
84     print('[INFO] Получить результат df_name_title в виде DataFrame с MutIndex')
85     time_start = datetime.datetime.now()
86     df_name_title = pd.DataFrame(list_name_title, columns=['name', 'rols'])
87     list_name_title = []
88     df_name_title.sort_values(by='name', inplace=True)
89     df_name_title.reset_index(drop=True, inplace=True)
90     use_sec = (datetime.datetime.now() - time_start).seconds
91     print('[INFO] DataFrame Конец преобразования, во времени: ', use_sec, ' s (', use_sec / 60, ')
92 min')
93
94     # print('>>' * 50)
95     # print('[INFO] 开始序列化 (备份) df_name_title')
96     # time_start = datetime.datetime.now()
97     # f = open(dump_name_title, 'wb')
98     # pickle.dump(obj=df_name_title, file=f)
99     # f.close()
100     # use_sec = (datetime.datetime.now() - time_start).seconds
101     # print('[INFO] 序列化 (备份) 结束, 用时: ', use_sec, ' s (', use_sec / 60, ') min')
102
103     """Объединить дубликаты имен, чтобы сделать их уникальными. добавить информацию об
104 одном и том же актере в rols"""
105     print('>>' * 50)
106     print('[INFO] Начните объединять дубликаты имен')
107     time_start = datetime.datetime.now()

```

```

105     name = None
106     tmp_df_rols = []
107     res_list_name_rols = []
108     for i in range(df_name_title.shape[0]):
109
110         """Если это один и тот же человек, работы объединяются в DataFrame"""
111         if name == df_name_title.loc[i]['name']:
112             tmp_df_rols = pd.concat([tmp_df_rols, df_name_title.loc[i]['rols']])
113
114         else:
115             """Приступая к следующему человеку
116             сначала записать информацию о предыдущем человеке в новый список, затем сбросить
117             имя и tmp_df_rols на содержимое текущего ряда
118             """
119             if name is not None:
120                 res_list_name_rols.append([name, tmp_df_rols])
121
122                 # if name == '$haniqua':
123                 #     print(tmp_df_rols)
124
125                 name = df_name_title.loc[i]['name']
126                 tmp_df_rols = df_name_title.loc[i]['rols']
127
128             if i % 10000 == 0:
129                 use_sec = (datetime.datetime.now() - time_start).seconds
130                 print('[INFO] Обработано ', i, ' строка | ', (i / df_name_title.shape[0]) * 100, '% | Затраченное
131                 время: ', use_sec, ' s (',
132                 use_sec / 60, ') min')
133
134             res_list_name_rols = pd.DataFrame(data=res_list_name_rols, columns=['name', 'rols'])
135
136             # print('>>' * 50)
137             # print('[INFO] 开始序列化（备份）res_list_name_rols')
138             # time_start = datetime.datetime.now()
139             # f = open(dump_df_res, 'wb')
140             # pickle.dump(obj=res_list_name_rols, file=f)
141             # f.close()
142             # use_sec = (datetime.datetime.now() - time_start).seconds
143             # print('[INFO] 序列化（备份）res_list_name_rols 结束，用时: ', use_sec, ' s (', use_sec / 60, ')
144             # min')
145
146             return res_list_name_rols
147
148 if __name__ == '__main__':
149     print('>>' * 50)
150     print('[INFO] Начало реализации')
151
152     #####
153     #####
154
155     # 演员信息表
156     print('>>' * 50)

```

```

153     print('[INFO] Начните читать `name.basics.tsv`')
154     time_start = datetime.datetime.now()
155     df_name_info = pd.read_csv(
156         filepath_or_buffer='/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt'
157         '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 семестр/СУБД/资
158         料/DataSet/name.basics.tsv',
159         header=0,
160         sep='\t'
161     )
162     df_name_info = df_name_info.iloc[:, :-1]
163     df_name_info.columns = ['nconst', 'name', 'birthYear', 'deathYear', 'profession']
164     time_end = datetime.datetime.now()
165     print('[INFO] Конец чтения `name.basics.tsv`, время: ', (time_end - time_start).seconds, ' s')
166
167     #####
168
169     # Обработка отсутствующих значений как None для простого преобразования в JSON
170     print('>>' * 50)
171     print('[INFO] Начните обрабатывать отсутствующие значения как None и удалите
172     дублирующиеся значения для легкого преобразования в JSON')
173     time_start = datetime.datetime.now()
174     df_name_info.replace(to_replace=['\\N', np.nan], value=None, inplace=True)
175     df_name_info.drop_duplicates(subset='name', keep='first', inplace=True)
176     time_end = datetime.datetime.now()
177     print('[INFO] Окончание обработки отсутствующих значений и удаления дублирующих
178     значений, время: ', (time_end - time_start).seconds, ' s')
179
180     #####
181
182     print('>>' * 50)
183     print('[INFO] Начните чтение файла, '/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt'
184     '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 '
185     'семестр/СУБД/资料/DataSet/data_actors.list.txt')
186     time_start = datetime.datetime.now()
187     # source_data = pd.read_csv(
188     #     filepath_or_buffer='/Users/fox/Library/CloudStorage/OneDrive-
189     #     PetertheGreatSt.PetersburgPolytechnicalUniversity/СПБПУ/3 '
190     #     'курс/6 семестр/СУБД/资料/DataSet/data_actresses.list.txt',
191     #     header=0,
192     #     sep='\t'
193     # )
194
195     # Список мужчин-актеров DataFrame
196     source_data = pd.read_csv(

```

```

194     filepath_or_buffer='/Users/fox/Library/CloudStorage/OneDrive-
PetertheGreatSt.PetersburgPolytechnicalUniversity/СПБПУ/3 '
195     'курс/6 семестр/СУБД/资料/DataSet/data_actors.list.txt',
196     header=0,
197     sep='\t'
198 )
199 time_end = datetime.datetime.now()
200 print('[INFO] Конец чтения файла, время: ', (time_end - time_start).seconds, ' s')
201
202
203
204 print('>>' * 50)
205 print('[INFO] Начните предварительную обработку данных')
206 time_start = datetime.datetime.now()
207 source_data.columns = ['name', 't1', 't2', 't3']
208 source_data['t1'].fillna(value="", inplace=True)
209 source_data['t2'].fillna(value="", inplace=True)
210 source_data['t3'].fillna(value="", inplace=True)
211
212 movie_list = source_data['t1'] + source_data['t2'] + source_data['t3']
213 source_data = pd.concat([source_data['name'], movie_list], axis=1)
214 source_data.columns = ['name', 'title_mix']
215
216 # source_data = source_data.iloc[:1000, :]
217
218 time_end = datetime.datetime.now()
219 print('[INFO] Окончание предварительной обработки данных, время: ', (time_end -
time_start).seconds, ' s')
220
221
222 start_index = 0
223 step_index = 1000000
224 end_index = start_index + step_index
225 times = 31
226
227
228
229 for i in range(1, times):
230     print('\n\n')
231     print('>>' * 20, ' Запустите цикл ', i, '<<' * 20)
232
233     tmp_source_data = source_data.iloc[start_index:end_index, :]
234
235
236     #####
237     #####
238     # # Список актрис df
239     # print('[INFO] main -> Начать обработку `data_actresses.list.txt`)
240     # df_actresses = normal_actor(
241     #     source_data=tmp_source_data,
242     #     dump_name_title='./result/dump_df_actresses_name_title.bits',
243     #     dump_df_res='./result/dump_df_actresses.bits'
244     # )

```

```

243     # tmp_source_data = None
244     # print('[INFO] main <- Окончание обработки `data_actresses.list.txt`)
245
246     #####
247     #####
248
249     #####
250     #####
251     # # 男演员列表df
252     # print('[INFO] Start handle `data_actors.list.txt`)
253     # df_actors = normal_actor(
254     #     file_path='/Users/fox/Library/CloudStorage/OneDrive-
PetertheGreatSt.PetersburgPolytechnicalUniversity/СПБПУ/3 '
255     #     'курс/6 семестр/СУБД/资料/DataSet/data_actors.list.txt'
256     # )
257     print('[INFO] main -> Начать обработку `data_actors.list.txt`)
258     df_actresses = normal_actor(
259     source_data=tmp_source_data,
260     dump_name_title='./result/dump_df_actors_name_title.bits',
261     dump_df_res='./result/dump_df_actors.bits'
262     )
263     tmp_source_data = None
264     print('[INFO] main <- Окончание обработки `data_actresses.list.txt`)
265
266     #####
267     #####
268
269     #####
270     #####
271
272     print('>>' * 50)
273     print('[INFO] Начните объединение двух больших таблиц для обработки в качестве
конечного результата')
274     time_start = datetime.datetime.now()
275     df_all = pd.merge(left=df_name_info,
276                       right=df_actresses,
277                       how='inner',
278                       on='name')
279     time_end = datetime.datetime.now()
280     print('[INFO] Конец слияния, во времени: ', (time_end - time_start).seconds, ' s')
281
282     df_actresses = None
283     print('[INFO] Освобождение памяти')
284
285     #####
286     #####
287
288     print('>>' * 50)
289     print('[INFO] Начните сериализацию (резервное копирование) конечного результата
слияния')
290     time_start = datetime.datetime.now()

```



```

283     dump_df_res = '/Users/fox/Library/CloudStorage/OneDrive-
PetertheGreatSt.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 семестр/СУБД/资
料/DataSet/result_dump_actors/dump_actors_' + str(i) + '.bits'
284     f = open(dump_df_res, 'wb')
285     pickle.dump(obj=df_all, file=f)
286     f.close()
287     use_sec = (datetime.datetime.now() - time_start).seconds
288     print('[INFO] Окончание сериализации (резервного копирования) объединенного конечного
результата в срок: ', use_sec, ' s (' , use_sec / 60, ') min')
289
290
291
#####
#####
292     print('[INFO] Start write to JSON file')
293     df_all.to_json(
294         path_or_buf='/Users/fox/Library/CloudStorage/OneDrive-
PetertheGreatSt.PetersburgPolytechnicalUniversity'
295             '/СПБПУ/3 курс/6 семестр/СУБД/资料/DataSet/result_json_actors/df_final_actors_' +
str(i) + '.json',
296         orient='records',
297         lines=True)
298     print('[INFO] JSON Вписать полностью')
299     df_all = None
300     print('[INFO] Освобождение памяти')
301
#####
#####
302
303     start_index += step_index
304     end_index += step_index
305
306     pass
307

```

2. Десериализовать все полученные малые наборы данных и сшить их вместе в большой набор данных, а также удалить дублирующиеся значения

```

1  | # -----*----- coding: utf-8 -----*-----
2  | # @Time   : 2022/4/27 14:37
3  | # @Author : 冰糖雪狸 (NekoSilverfox)
4  | # @Project : IMDB 数据处理
5  | # @File   : 反序列化数据.py.py
6  | # @Software: PyCharm
7  | # @Github : https://github.com/NekoSilverFox
8  | # -----
9  | import pickle
10 | import datetime
11 | import pandas as pd
12 |
13 |

```

```

14 def concat_df(bits_file_path_header: str,
15               max_index: int,
16               path_result_bits_save: str) -> pd.DataFrame:
17     """
18     Десериализовать из многих сериализованных файлов и соединить их вместе
19     :param bits_file_path_header: Сериализация [заголовка] файла
20     :param max_index: Максимальный индекс заголовка документа
21     :param path_result_bits_save: Где сохранить десериализацию результата слияния
22     :return: DataFrame после десериализации
23     """
24     df_result = None
25
26     for i in range(1, max_index + 1):
27         print('-' * 50)
28         print('[INFO] Начните читать первый файл ', i)
29         time_start = datetime.datetime.now()
30         bits_file_path = bits_file_path_header + str(i) + '.bits'
31         f = open(bits_file_path, 'rb')
32         df_obj = pickle.load(file=f)
33         f.close()
34         time_end = datetime.datetime.now()
35         print('[INFO] Прочитайте конец файла ', i, ' | Затраченное время: ', (time_end -
time_start).seconds, ' s\n')
36
37         if i == 1:
38             df_result = df_obj
39             continue
40
41         print('[INFO] Начните сращивание файла ', i)
42         time_start = datetime.datetime.now()
43         df_result = pd.concat([df_result, df_obj])
44         time_end = datetime.datetime.now()
45         print('[INFO] Конец ', i, 'сращивания документов, Затраченное время: ', (time_end -
time_start).seconds, ' s\n')
46
47         # Завершение слияния, сохранение результатов с помощью сериализации
48         print('-' * 50)
49         print('[INFO] Завершение слияния, сохранение результатов с помощью сериализации')
50         time_start = datetime.datetime.now()
51         f = open(path_result_bits_save, 'wb')
52         pickle.dump(obj=df_result, file=f)
53         f.close()
54         time_end = datetime.datetime.now()
55         print('[INFO] Конец сериализации для сохранения результатов, во времени: ', (time_end -
time_start).seconds, ' s\n')
56
57     return df_result
58
59
60 def merge_duplicates(df_source: pd.DataFrame) -> pd.DataFrame:
61     """
62     Повторное дублирование конечного результата
63     :param df_source: Слияние мужских и женских актеров - это массив

```

```

64 :return: Объединенные и де-дублицированные DataFrame
65 """
66 df_source.sort_values(by='nconst', inplace=True)
67 df_source.reset_index(drop=True, inplace=True)
68
69 i_current = 0
70 i_next = i_current + 1
71 stop_index = df_source.shape[0]
72 while i_next <= stop_index:
73     while df_source.loc[i_current]['nconst'] == df_source.loc[i_next]['nconst']:
74         df_source.loc[i_current]['rols'] = pd.concat([df_source.loc[i_current]['rols'],
df_source.loc[i_next]['rols']])
75         df_source.drop(index=i_next, inplace=True)
76         i_next += 1
77
78     if i_next == stop_index:
79         df_source.reset_index(drop=True, inplace=True)
80         return df_source
81
82     print('[INFO] объединенный ', i_current, ' ряд | ', round(i_current / df_source.shape[0] * 100, 4),
'%%')
83     i_current = i_next
84     i_next += 1
85
86 df_source.reset_index(drop=True, inplace=True)
87 return df_source
88
89
90 if __name__ == '__main__':
91
92     #####
93     # Объединить всех мужчин-актеров (актеров)
94
95     #####
96
97     print('>>' * 50)
98     bits_file_path_header = '/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt' \
99         '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 '\
100         'семестр/СУБД/资料/DataSet/result_dump_actors/dump_actors_'
101
102     path_result_bits_save = '/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt' \
103         '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 '\
104         'семестр/СУБД/资料/DataSet/result_dump_actors/dump_actors_ALL.bits'
105
106     df_all_actors = concat_df(bits_file_path_header=bits_file_path_header,
107         max_index=20,
108         path_result_bits_save=path_result_bits_save)
109     print('[INFO] Слияние и сериализация вывода успешно! \n выход на: ', path_result_bits_save)
110

```

```

109 #####
110 # Объединить всех актрис (actresses)
111 #####
112 print('>>' * 50)
113 bits_file_path_header = '/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt' \
114     '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 '\
115     'семестр/СУБД/资料/DataSet/result_dump_actresses/dump_actresses_'
116
117 path_result_bits_save = '/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt' \
118     '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 '\
119     'семестр/СУБД/资料/DataSet/result_dump_actresses/dump_actresses_ALL.bits'
120
121 df_all_actresses = concat_df(bits_file_path_header=bits_file_path_header,
122     max_index=13,
123     path_result_bits_save=path_result_bits_save)
124 print('[INFO] Слияние и сериализация вывода успешно! \n выход на: ', path_result_bits_save)
125
126
127 #####
128 # Объединить всех мужчин-актеров (актеров) и женщин-актеров (актрис)
129 #####
130 print('>>' * 50)
131 print('[INFO] Объединить всех мужчин-актеров (актеров) и женщин-актеров (актрис)')
132 time_start = datetime.datetime.now()
133
134 df_result_all = pd.concat([df_all_actors, df_all_actresses])
135
136 time_end = datetime.datetime.now()
137 print('[INFO] Объедините всех мужчин-актеров (актеров) и женщин-актеров (актрис) до конца,
138 в.', (time_end - time_start).seconds, ' s\n')
139
140 #####
141 df_all_actors = None
142 df_all_actresses = None
143 print('[INFO] Освобождение памяти')
144 #####
145
146 #####

```

```

147 # 合并后的最终结果再次去重
148
#####
#####
149 print('>>' * 50)
150 print('[INFO] Конечный результат слияния снова дедуплицируется')
151 time_start = datetime.datetime.now()
152 df_result_all = merge_duplicates(df_source=df_result_all)
153 time_end = datetime.datetime.now()
154 print('[INFO] Конечный результат слияния снова дедуплицируется для завершения: ',
(time_end - time_start).seconds, ' s\n') # 9714 s
155
156
#####
#####
157 # Окончательный результат слияния затем сериализуется и сохраняется в виде JSON-файла
158
#####
#####
159 path_result_bits_save = '/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt' \
160                          '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 '\
161                          'семестр/СУБД/资料/DataSet/result_ALL/dump_ALL.bits'
162 print('>>' * 50)
163 print('[INFO] Конец слияния, сохранение [окончательного] результата с помощью
сериализации')
164 time_start = datetime.datetime.now()
165 f = open(path_result_bits_save, 'wb')
166 pickle.dump(obj=df_result_all, file=f)
167 f.close()
168 time_end = datetime.datetime.now()
169 print('[INFO] Конец сериализации для сохранения результатов, во времени: ', (time_end -
time_start).seconds, ' s\n')
170
171
172 path_result_json_save = '/Users/fox/Library/CloudStorage/OneDrive-PetertheGreatSt' \
173                          '.PetersburgPolytechnicalUniversity/СПБПУ/3 курс/6 '\
174                          'семестр/СУБД/资料/DataSet/result_ALL/dump_ALL.json'
175 print('>>' * 50)
176 print('[INFO] Сохраните [окончательный] результат в формате JSON')
177 time_start = datetime.datetime.now()
178 df_result_all.to_json(path_or_buf=path_result_json_save,
179                       orient='records',
180                       lines=True)
181 time_end = datetime.datetime.now()
182 print('[INFO] Сохраните [окончательный] результат в формате JSON в формате: ', (time_end -
time_start).seconds, ' s\n')
183
#####
#####
184 pass
185

```

3. Выполнение тестов на скорость, сериализация и построение результатов в виде изображений

```

1  # -----*----- coding: utf-8 -----*-----
2  # @Time   : 2022/4/28 14:22
3  # @Author : 冰糖雪狸 (NekoSilverfox)
4  # @Project : JSON 速度测试
5  # @File   : 速度测试.py
6  # @Software: PyCharm
7  # @Github : https://github.com/NekoSilverFox
8  # -----
9  import psycopg2 as pg
10 import matplotlib.pyplot as plt
11 import pandas as pd
12 import datetime
13 import pickle
14
15
16 def get_id_len_json_df():
17     """
18     Получение длины данных JSON, соответствующих каждому идентификатору в базе данных
19     :return:
20     """
21     # Получение длины данных JSON, соответствующих каждому идентификатору в базе данных
22     print('>>' * 50)
23     print('[INFO] Start connect database')
24     conn = pg.connect(database="db_imdb",
25                       user="postgres",
26                       password="postgres",
27                       host="localhost",
28                       port="5432")
29     cur = conn.cursor()
30     print('[INFO] Connect database successfully')
31
32     # Получить количество строк
33     cur.execute("SELECT COUNT(*) FROM tb_json;")
34     count_row = cur.fetchall()[0][0]
35     print('количество строк: ', count_row)
36
37     # DataFrame для статистики
38     df_counter = pd.DataFrame([[0, 0]], columns=['id', 'len_json'])
39
40     # Выполните запрос и засекайте время
41     for i in range(1, count_row):
42         print('[INFO] тестирует строку ' + str(i) + ' строка | ' + str(round(i / count_row * 100, 4)) + '%')
43
44         comm_sql = 'SELECT imdata FROM tb_json WHERE iddata=' + str(i) + ';'
45         cur.execute(comm_sql)
46         len_row_json = len(str(cur.fetchall()[0])) # длина JSON(B)
47
48         df_tmp = pd.DataFrame([[i, len_row_json]], columns=['id', 'len_json'])
49         df_counter = pd.concat([df_counter, df_tmp])
50
51     conn.close()

```

```

52
53 df_counter = df_counter.iloc[1:, :]
54 df_counter.sort_values(by='len_json', inplace=True)
55
56 print('>>' * 50)
57 print('[INFO] Конец слияния, сохранение [окончательного] результата с помощью
сериализации')
58 time_start = datetime.datetime.now()
59 f = open('./result/df_id_json_len.bits', 'wb')
60 pickle.dump(obj=df_counter, file=f)
61 f.close()
62 time_end = datetime.datetime.now()
63 print('[INFO] Конец сериализации для сохранения результатов, во времени: ', (time_end -
time_start).seconds, ' \n')
64
65
66 def json_by_id_every_col_test():
67     # Получение длины данных JSON, соответствующих каждому идентификатору в базе данных
68     print('>>' * 50)
69     print('[INFO] Start connect database')
70     conn = pg.connect(database="db_imdb",
71                       user="postgres",
72                       password="postgres",
73                       host="localhost",
74                       port="5432")
75     cur = conn.cursor()
76     print('[INFO] Connect database successfully')
77
78     # Получить количество строк
79     cur.execute("SELECT COUNT(*) FROM tb_json;")
80     count_row = cur.fetchall()[0][0]
81     print('Количество строк: ', count_row)
82
83     # DataFrame для статистики
84     df_counter = pd.DataFrame([[0, 0, 0, 0, 0, 0, 0]],
85                               columns=['len_row', 'full_ms', 'nconst_ms', 'name_ms', 'birthYear_ms',
'profession_ms', 'rols_ms'])
86
87     # Выполните запрос и засекайте время
88     for i in range(1, count_row):
89         print('[INFO] тестирует строку ' + str(i) + ' строка | ' + str(round(i / count_row * 100, 4)) + '%')
90
91         # Все поля во всей строке данных
92         comm_sql = 'SELECT imdata FROM tb_json WHERE iddata=' + str(i) + ';'
93         start_time = datetime.datetime.now()
94         cur.execute(comm_sql)
95         end_time = datetime.datetime.now()
96         full_use_time_ms = (end_time - start_time).microseconds
97         row = cur.fetchall()[0]
98         len_row_json = len(str(row)) # длина JSON(B)
99
100         # Поля во всей строке данных - nconst
101         comm_sql = "SELECT imdata->>'nconst' FROM tb_json WHERE iddata=" + str(i) + ';'

```

```

102 start_time = datetime.datetime.now()
103 cur.execute(comm_sql)
104 end_time = datetime.datetime.now()
105 nconst_use_time_ms = (end_time - start_time).microseconds
106
107 # Поля во всей строке данных - name
108 comm_sql = "SELECT imdata->>'name' FROM tb_json WHERE iddata=" + str(i) + ';'
109 start_time = datetime.datetime.now()
110 cur.execute(comm_sql)
111 end_time = datetime.datetime.now()
112 name_use_time_ms = (end_time - start_time).microseconds
113
114 # Поля во всей строке данных - birthYear
115 comm_sql = "SELECT imdata->>'birthYear' FROM tb_json WHERE iddata=" + str(i) + ';'
116 start_time = datetime.datetime.now()
117 cur.execute(comm_sql)
118 end_time = datetime.datetime.now()
119 birthYear_use_time_ms = (end_time - start_time).microseconds
120
121 # Поля во всей строке данных - profession
122 comm_sql = "SELECT imdata->>'profession' FROM tb_json WHERE iddata=" + str(i) + ';'
123 start_time = datetime.datetime.now()
124 cur.execute(comm_sql)
125 end_time = datetime.datetime.now()
126 profession_use_time_ms = (end_time - start_time).microseconds
127
128 # Поля во всей строке данных - rols
129 comm_sql = "SELECT imdata->>'rols' FROM tb_json WHERE iddata=" + str(i) + ';'
130 start_time = datetime.datetime.now()
131 cur.execute(comm_sql)
132 end_time = datetime.datetime.now()
133 rols_use_time_ms = (end_time - start_time).microseconds
134
135 df_tmp = pd.DataFrame([[len_row_json, full_use_time_ms, nconst_use_time_ms,
name_use_time_ms, birthYear_use_time_ms, profession_use_time_ms, rols_use_time_ms]],
136                        columns=['len_row', 'full_ms', 'nconst_ms', 'name_ms', 'birthYear_ms',
'profession_ms', 'rols_ms'])
137 df_counter = pd.concat([df_counter, df_tmp])
138
139 conn.close()
140
141 df_counter = df_counter.iloc[1:, :]
142 df_counter.sort_values(by='len_row', inplace=True)
143
144 print('>>' * 50)
145 print('[INFO] Конец слияния, сохранение [окончательного] результата с помощью
сериализации')
146 time_start = datetime.datetime.now()
147 f = open('./result/json/res_id_ix_every_col.bits', 'wb')
148 pickle.dump(obj=df_counter, file=f)
149 f.close()
150 time_end = datetime.datetime.now()

```



```

151     print('[INFO] Конец сериализации для сохранения результатов, во времени: ', (time_end -
time_start).seconds, ' s\n')
152
153     # print(df_counter)
154
155     """Построение графиков результатов"""
156     plt.figure(figsize=(20, 10), dpi=100)
157     plt.scatter(x=df_counter['len_row'].values,
158                y=df_counter['full_ms'].values,
159                label='full row')
160     plt.scatter(x=df_counter['len_row'].values,
161                y=df_counter['nconst_ms'].values,
162                label='nconst')
163     plt.scatter(x=df_counter['len_row'].values,
164                y=df_counter['name_ms'].values,
165                label='name')
166     plt.scatter(x=df_counter['len_row'].values,
167                y=df_counter['birthYear_ms'].values,
168                label='birthYear')
169     plt.scatter(x=df_counter['len_row'].values,
170                y=df_counter['profession_ms'].values,
171                label='profession')
172     plt.scatter(x=df_counter['len_row'].values,
173                y=df_counter['rols_ms'].values,
174                label='rols')
175     plt.legend()
176     plt.title('Query by key `ID` in tb_json')
177     plt.xlabel('Length of JSON')
178     plt.ylabel('Query time (milliseconds)')
179     plt.savefig('./result/json/res_id_ix_every_col.png')
180     # plt.show()
181
182
183     def jsonb_by_id_every_col_test():
184         # Получение длины данных JSON, соответствующих каждому идентификатору в базе данных
185         print('>>' * 50)
186         print('[INFO] Start connect database')
187         conn = pg.connect(database="db_imdb",
188                           user="postgres",
189                           password="postgres",
190                           host="localhost",
191                           port="5432")
192         cur = conn.cursor()
193         print('[INFO] Connect database successfully')
194
195         # Получить количество строк
196         cur.execute("SELECT COUNT(*) FROM tb_jsonb;")
197         count_row = cur.fetchall()[0][0]
198         print('Количество строк: ', count_row)
199
200         # DataFrame для статистики
201         df_counter = pd.DataFrame([[0, 0, 0, 0, 0, 0, 0]],

```

```

202         columns=['len_row', 'full_ms', 'nconst_ms', 'name_ms', 'birthYear_ms',
'profession_ms', 'rols_ms'])
203
204     # Выполните запрос и засекайте время
205     for i in range(1, count_row):
206         print('[INFO] тестирует строку ' + str(i) + ' строка | ' + str(round(i / count_row * 100, 4)) + '%')
207
208         # Все поля во всей строке данных
209         comm_sql = 'SELECT imdata FROM tb_jsonb WHERE iddata=' + str(i) + ';'
210         start_time = datetime.datetime.now()
211         cur.execute(comm_sql)
212         end_time = datetime.datetime.now()
213         full_use_time_ms = (end_time - start_time).microseconds
214         row = cur.fetchall()[0]
215         len_row_jsonb = len(str(row)) # длина JSON(B)
216
217         # Поля во всей строке данных - nconst
218         comm_sql = "SELECT imdata->>'nconst' FROM tb_jsonb WHERE iddata=" + str(i) + ';'
219         start_time = datetime.datetime.now()
220         cur.execute(comm_sql)
221         end_time = datetime.datetime.now()
222         nconst_use_time_ms = (end_time - start_time).microseconds
223
224         # Поля во всей строке данных - name
225         comm_sql = "SELECT imdata->>'name' FROM tb_jsonb WHERE iddata=" + str(i) + ';'
226         start_time = datetime.datetime.now()
227         cur.execute(comm_sql)
228         end_time = datetime.datetime.now()
229         name_use_time_ms = (end_time - start_time).microseconds
230
231         # Поля во всей строке данных - birthYear
232         comm_sql = "SELECT imdata->>'birthYear' FROM tb_jsonb WHERE iddata=" + str(i) + ';'
233         start_time = datetime.datetime.now()
234         cur.execute(comm_sql)
235         end_time = datetime.datetime.now()
236         birthYear_use_time_ms = (end_time - start_time).microseconds
237
238         # Поля во всей строке данных - profession
239         comm_sql = "SELECT imdata->>'profession' FROM tb_jsonb WHERE iddata=" + str(i) + ';'
240         start_time = datetime.datetime.now()
241         cur.execute(comm_sql)
242         end_time = datetime.datetime.now()
243         profession_use_time_ms = (end_time - start_time).microseconds
244
245         # Поля во всей строке данных - rols
246         comm_sql = "SELECT imdata->>'rols' FROM tb_jsonb WHERE iddata=" + str(i) + ';'
247         start_time = datetime.datetime.now()
248         cur.execute(comm_sql)
249         end_time = datetime.datetime.now()
250         rols_use_time_ms = (end_time - start_time).microseconds
251
252         df_tmp = pd.DataFrame([[len_row_jsonb, full_use_time_ms, nconst_use_time_ms,
name_use_time_ms, birthYear_use_time_ms, profession_use_time_ms, rols_use_time_ms]],

```

```

253         columns=['len_row', 'full_ms', 'nconst_ms', 'name_ms', 'birthYear_ms',
'profession_ms', 'rols_ms'])
254     df_counter = pd.concat([df_counter, df_tmp])
255
256     conn.close()
257
258     df_counter = df_counter.iloc[1:, :]
259     df_counter.sort_values(by='len_row', inplace=True)
260
261     print('>>' * 50)
262     print('[INFO] Конец слияния, сохранение [окончательного] результата с помощью
сериализации')
263     time_start = datetime.datetime.now()
264     f = open('./result/jsonb/res_id_ix_every_col.bits', 'wb')
265     pickle.dump(obj=df_counter, file=f)
266     f.close()
267     time_end = datetime.datetime.now()
268     print('[INFO] Конец сериализации для сохранения результатов, во времени: ', (time_end -
time_start).seconds, ' \n')
269
270     # print(df_counter)
271
272     """Построение графиков результатов"""
273     plt.figure(figsize=(20, 10), dpi=100)
274     plt.scatter(x=df_counter['len_row'].values,
275               y=df_counter['full_ms'].values,
276               label='full row')
277     plt.scatter(x=df_counter['len_row'].values,
278               y=df_counter['nconst_ms'].values,
279               label='nconst')
280     plt.scatter(x=df_counter['len_row'].values,
281               y=df_counter['name_ms'].values,
282               label='name')
283     plt.scatter(x=df_counter['len_row'].values,
284               y=df_counter['birthYear_ms'].values,
285               label='birthYear')
286     plt.scatter(x=df_counter['len_row'].values,
287               y=df_counter['profession_ms'].values,
288               label='profession')
289     plt.scatter(x=df_counter['len_row'].values,
290               y=df_counter['rols_ms'].values,
291               label='rols')
292     plt.legend()
293     plt.title('Query by key `ID` in tb_jsonb')
294     plt.xlabel('Length of JSONB')
295     plt.ylabel('Query time (milliseconds)')
296     plt.savefig('./result/jsonb/res_id_ix_every_col.png')
297     # plt.show()
298
299
300 def json_update_by_id_every_col_test():
301     # Получение длины данных JSON, соответствующих каждому идентификатору в базе данных
302     print('>>' * 50)

```

```

303 print('[INFO] Start connect database')
304 conn = pg.connect(database="db_imdb",
305                   user="postgres",
306                   password="postgres",
307                   host="localhost",
308                   port="5432")
309 cur = conn.cursor()
310 print('[INFO] Connect database successfully')
311
312 # Получить количество строк
313 cur.execute("SELECT COUNT(*) FROM tb_json;")
314 count_row = cur.fetchall()[0][0]
315 print('Количество строк: ', count_row)
316
317 # DataFrame для статистики
318 df_counter = pd.DataFrame([[0, 0, 0, 0, 0]],
319                           columns=['len_row', 'nconst_ms', 'name_ms', 'birthYear_ms', 'rols_ms'])
320
321 # Выполните запрос и засекуте время
322 for i in range(1, count_row):
323     print("\n[INFO] тестирует строку ' + str(i) + ' строка | ' + str(round(i / count_row * 100, 4)) + '%')
324
325     try:
326         # Все поля во всей строке данных
327         comm_sql = 'SELECT imdata FROM tb_json WHERE iddata=' + str(i) + ';'
328         cur.execute(comm_sql)
329         len_row_json = len(str(cur.fetchall()[0])) # длина JSON(B)
330
331         cur.execute('BEGIN;')
332
333         print("\tBEGIN;")
334
335         # Поля во всей строке данных - nconst
336         comm_sql = "UPDATE tb_json SET imdata=jsonb_set(imdata::jsonb, '{nconst}',
337 \"tt0000009\"::jsonb) WHERE iddata=" + str(i) + ';'
338         start_time = datetime.datetime.now()
339         cur.execute(comm_sql)
340         end_time = datetime.datetime.now()
341         nconst_use_time_ms = (end_time - start_time).microseconds
342         print("\tnconst Конец испытания, время: ', nconst_use_time_ms, ' ms')
343
344         # Поля во всей строке данных - name
345         comm_sql = "UPDATE tb_json SET imdata=jsonb_set(imdata::jsonb, '{name}',
346 \"tt_name\"::jsonb) WHERE iddata=" + str(i) + ';'
347         start_time = datetime.datetime.now()
348         cur.execute(comm_sql)
349         end_time = datetime.datetime.now()
350         name_use_time_ms = (end_time - start_time).microseconds
351         print("\tname Конец испытания, время: ', name_use_time_ms, ' ms')
352
353         # Поля во всей строке данных - birthYear
354         comm_sql = "UPDATE tb_json SET imdata=jsonb_set(imdata::jsonb, '{birthYear}',
355 \"2222\"::jsonb) WHERE iddata=" + str(i) + ';'

```

```

353     start_time = datetime.datetime.now()
354     cur.execute(comm_sql)
355     end_time = datetime.datetime.now()
356     birthYear_use_time_ms = (end_time - start_time).microseconds
357     print("\tbirthYear Конец испытания, время: ', birthYear_use_time_ms, ' ms')
358
359     # Поля во всей строке данных - rols
360     comm_sql = 'UPDATE tb_json SET imdata=jsonb_set(imdata::jsonb, \'{"year": 2000,
"title": "t_title", "series name": "t_series", "character name": "t_character_name"}\'::jsonb) WHERE
iddata=' + str(i) + ';'
361     start_time = datetime.datetime.now()
362     cur.execute(comm_sql)
363     end_time = datetime.datetime.now()
364     rols_use_time_ms = (end_time - start_time).microseconds
365     print("\trols Конец испытания, время: ', rols_use_time_ms, ' ms')
366
367     except:
368         cur.execute('ROLLBACK;')
369         continue
370
371     cur.execute('ROLLBACK;')
372     df_tmp = pd.DataFrame([[len_row_json, nconst_use_time_ms, name_use_time_ms,
birthYear_use_time_ms, rols_use_time_ms]],
373                           columns=['len_row', 'nconst_ms', 'name_ms', 'birthYear_ms', 'rols_ms'])
374     df_counter = pd.concat([df_counter, df_tmp])
375
376     conn.close()
377
378     df_counter = df_counter.iloc[1:, :]
379     df_counter.sort_values(by='len_row', inplace=True)
380
381     print('>>' * 50)
382     print('[INFO] Конец слияния, сохранение [окончательного] результата с помощью
сериализации')
383     time_start = datetime.datetime.now()
384     f = open('./result/json/res_update_id_ix_every_col.bits', 'wb')
385     pickle.dump(obj=df_counter, file=f)
386     f.close()
387     time_end = datetime.datetime.now()
388     print('[INFO] Конец сериализации для сохранения результатов, во времени: ', (time_end -
time_start).seconds, ' s\n')
389
390     # print(df_counter)
391
392     """Построение графиков результатов"""
393     plt.figure(figsize=(20, 10), dpi=100)
394     plt.scatter(x=df_counter['len_row'].values,
395               y=df_counter['nconst_ms'].values,
396               label='nconst')
397     plt.scatter(x=df_counter['len_row'].values,
398               y=df_counter['name_ms'].values,
399               label='name')
400     plt.scatter(x=df_counter['len_row'].values,

```

```

401     y=df_counter['birthYear_ms'].values,
402     label='birthYear')
403 plt.scatter(x=df_counter['len_row'].values,
404             y=df_counter['rols_ms'].values,
405             label='rols')
406 plt.legend()
407 plt.title('Test UPDATE, query by key `ID` in tb_json')
408 plt.xlabel('Length of JSON')
409 plt.ylabel('Query time (milliseconds)')
410 plt.savefig('./result/json/res_update_id_ix_every_col.png')
411 # plt.show()
412
413
414 def jsonb_update_by_id_every_col_test():
415     # Получение длины данных JSON, соответствующих каждому идентификатору в базе данных
416     print('>>' * 50)
417     print('[INFO] Start connect database')
418     conn = pg.connect(database="db_imdb",
419                       user="postgres",
420                       password="postgres",
421                       host="localhost",
422                       port="5432")
423     cur = conn.cursor()
424     print('[INFO] Connect database successfully')
425
426     # Получить количество строк
427     cur.execute("SELECT COUNT(*) FROM tb_jsonb;")
428     count_row = cur.fetchall()[0][0]
429     print('Количество строк: ', count_row)
430
431     # DataFrame для статистики
432     df_counter = pd.DataFrame([[0, 0, 0, 0, 0]],
433                               columns=['len_row', 'nconst_ms', 'name_ms', 'birthYear_ms', 'rols_ms'])
434
435     # Выполните запрос и засекайте время
436     for i in range(1, count_row):
437         print("\n[INFO] тестирует строку " + str(i) + ' строка | ' + str(round(i / count_row * 100, 4)) + '%')
438
439         try:
440             # Все поля во всей строке данных
441             comm_sql = 'SELECT imdata FROM tb_jsonb WHERE iddata=' + str(i) + ';'
442             cur.execute(comm_sql)
443             len_row_jsonb = len(str(cur.fetchall()[0])) # длина JSON(B)
444
445             cur.execute('BEGIN;')
446
447             print("\tBEGIN;")
448
449             # Поля во всей строке данных - nconst
450             comm_sql = "UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{nconst}',
451 '\tt0000009\':"::jsonb) WHERE iddata=" + str(i) + ';'
451             # comm_sql = "UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{nconst}',
452 '\tt0000009\':"::jsonb) WHERE iddata=" + str(i) + ';'

```

```

452     # print(comm_sql)
453     start_time = datetime.datetime.now()
454     cur.execute(comm_sql)
455     end_time = datetime.datetime.now()
456     nconst_use_time_ms = (end_time - start_time).microseconds
457     print("\tnconst Конец испытания, время: ', nconst_use_time_ms, ' ms')
458
459     # Поля во всей строке данных - name
460     comm_sql = "UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{name}',
'\tt_name\':"::jsonb) WHERE iddata=" + str(i) + ';'
461     start_time = datetime.datetime.now()
462     cur.execute(comm_sql)
463     end_time = datetime.datetime.now()
464     name_use_time_ms = (end_time - start_time).microseconds
465     print("\tname Конец испытания, время: ', name_use_time_ms, ' ms')
466
467     # Поля во всей строке данных - birthYear
468     comm_sql = "UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{birthYear}',
'\t2222\':"::jsonb) WHERE iddata=" + str(i) + ';'
469     start_time = datetime.datetime.now()
470     cur.execute(comm_sql)
471     end_time = datetime.datetime.now()
472     birthYear_use_time_ms = (end_time - start_time).microseconds
473     print("\tbirthYear Конец испытания, время: ', birthYear_use_time_ms, ' ms')
474
475     # Поля во всей строке данных - rols
476     comm_sql = 'UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, \'{"year":
2000, "title": "t_title", "series name": "t_series", "character name": "t_character_name"}\'::jsonb)
WHERE iddata=' + str(i) + ';'
477     start_time = datetime.datetime.now()
478     cur.execute(comm_sql)
479     end_time = datetime.datetime.now()
480     rols_use_time_ms = (end_time - start_time).microseconds
481     print("\trols Конец испытания, время: ', rols_use_time_ms, ' ms')
482
483     except:
484         cur.execute('ROLLBACK;')
485         continue
486
487     cur.execute('ROLLBACK;')
488     df_tmp = pd.DataFrame([[len_row_jsonb, nconst_use_time_ms, name_use_time_ms,
birthYear_use_time_ms, rols_use_time_ms]],
489                           columns=['len_row', 'nconst_ms', 'name_ms', 'birthYear_ms', 'rols_ms'])
490     df_counter = pd.concat([df_counter, df_tmp])
491
492     conn.close()
493
494     df_counter = df_counter.iloc[1:, :]
495     df_counter.sort_values(by='len_row', inplace=True)
496
497     print('>' * 50)
498     print('[INFO] Конец слияния, сохранение [окончательного] результата с помощью
сериализации')

```

```

499     time_start = datetime.datetime.now()
500     f = open('./result/jsonb/res_update_id_ix_every_col.bits', 'wb')
501     pickle.dump(obj=df_counter, file=f)
502     f.close()
503     time_end = datetime.datetime.now()
504     print('[INFO] Конец сериализации для сохранения результатов, во времени: ', (time_end -
time_start).seconds, ' s\n')
505
506     # print(df_counter)
507
508     """Построение графиков результатов"""
509     plt.figure(figsize=(20, 10), dpi=100)
510     plt.scatter(x=df_counter['len_row'].values,
511                y=df_counter['nconst_ms'].values,
512                label='nconst')
513     plt.scatter(x=df_counter['len_row'].values,
514                y=df_counter['name_ms'].values,
515                label='name')
516     plt.scatter(x=df_counter['len_row'].values,
517                y=df_counter['birthYear_ms'].values,
518                label='birthYear')
519     plt.scatter(x=df_counter['len_row'].values,
520                y=df_counter['rols_ms'].values,
521                label='rols')
522     plt.legend()
523     plt.title('Test UPDATE, query by key `ID` in tb_jsonb')
524     plt.xlabel('Length of JSONB')
525     plt.ylabel('Query time (milliseconds)')
526     plt.savefig('./result/jsonb/res_update_id_ix_every_col.png')
527     # plt.show()
528
529
530 if __name__ == '__main__':
531     # get_id_len_json_df()
532
533     json_by_id_every_col_test()
534
535     jsonb_by_id_every_col_test()
536
537     json_update_by_id_every_col_test()
538
539     jsonb_update_by_id_every_col_test()
540

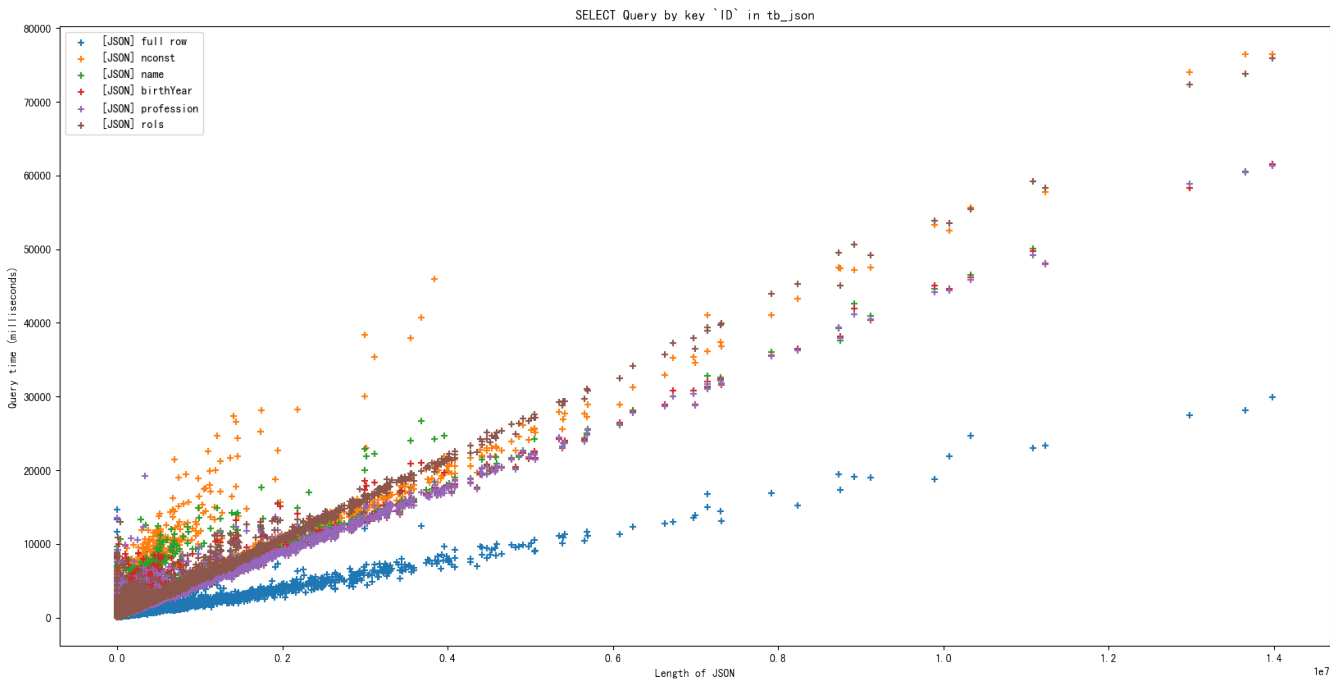
```


1.2.2.1.1 Результат

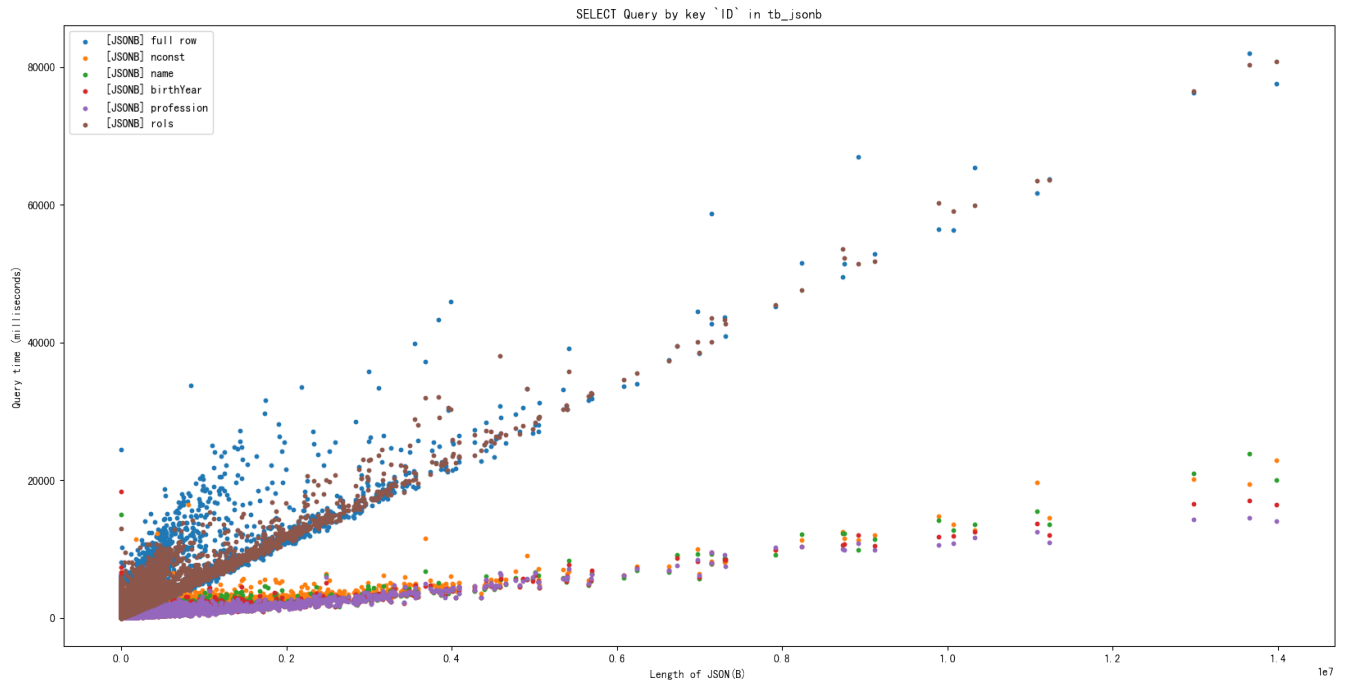
Сравнение времени, затрачиваемого JSON и JSONB при вставке данных в базу данных:

Типа	Время
JSON	49s
JSONB	79.32s

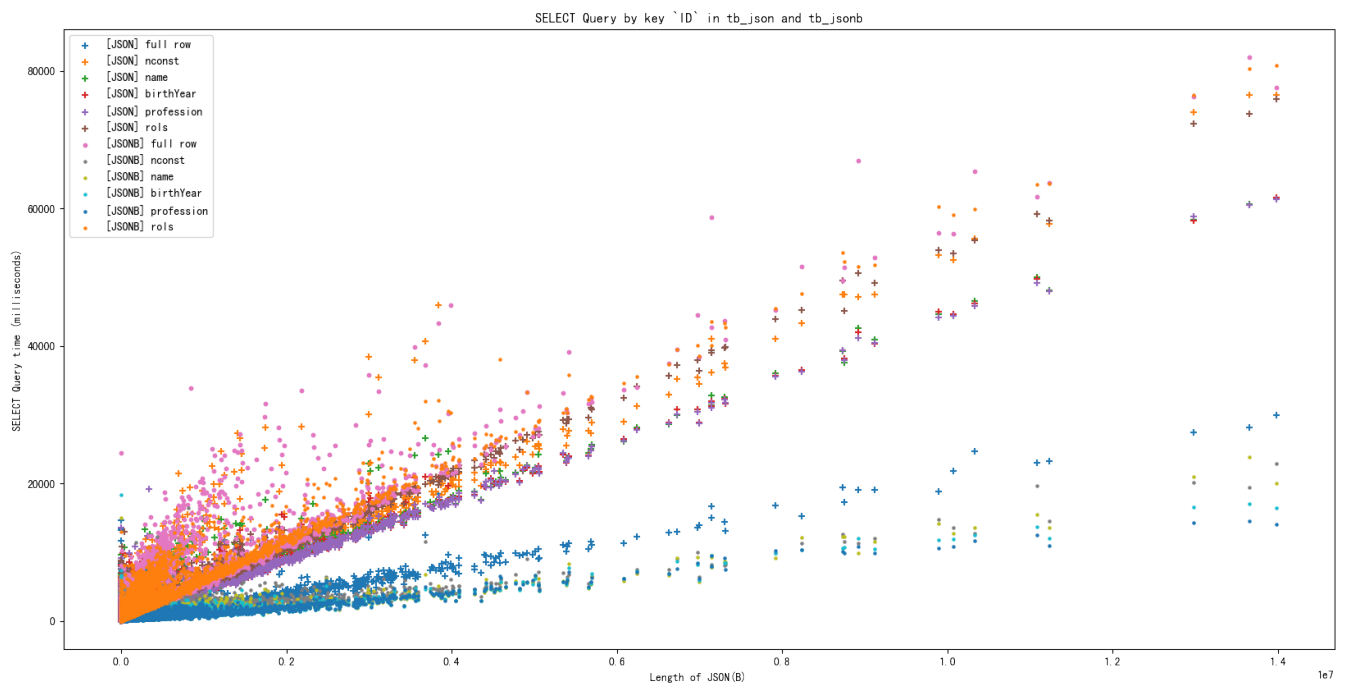
Сравнение времени SELECT для каждого поля JSON в таблице tb_json :



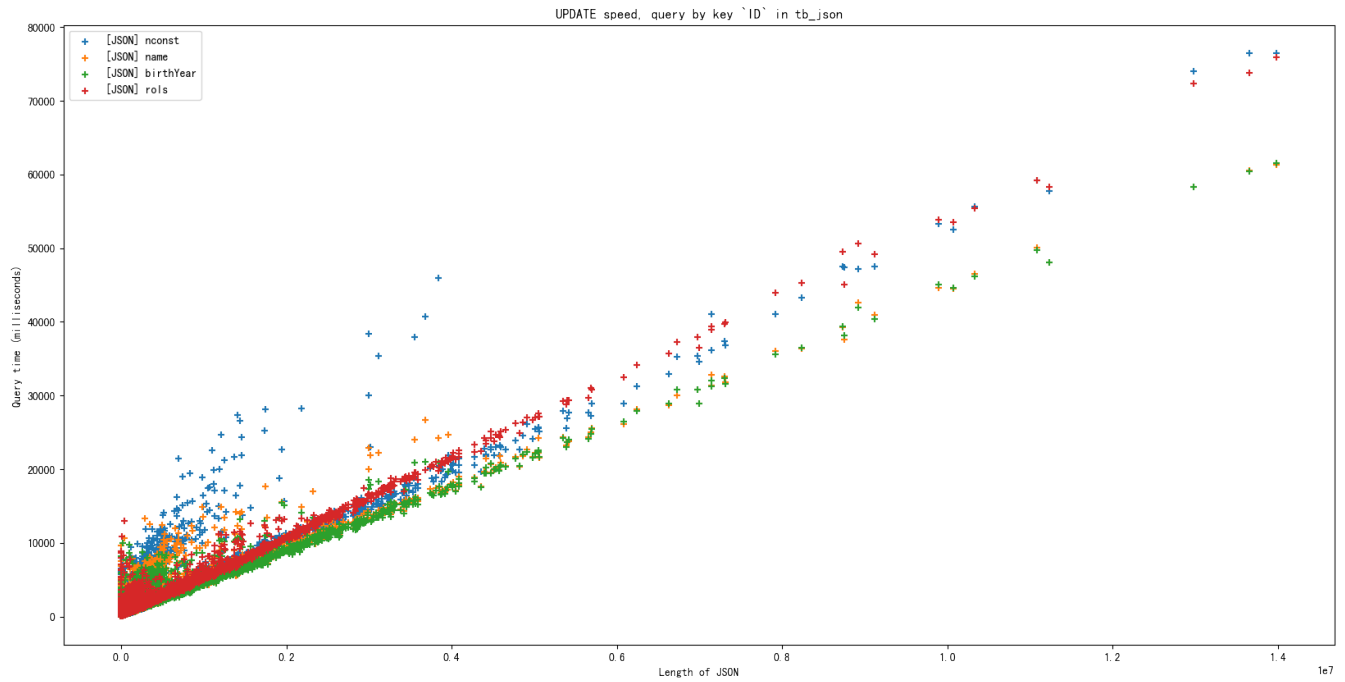
Сравнение времени SELECT для каждого поля JSONB в таблице tb_jsonb :



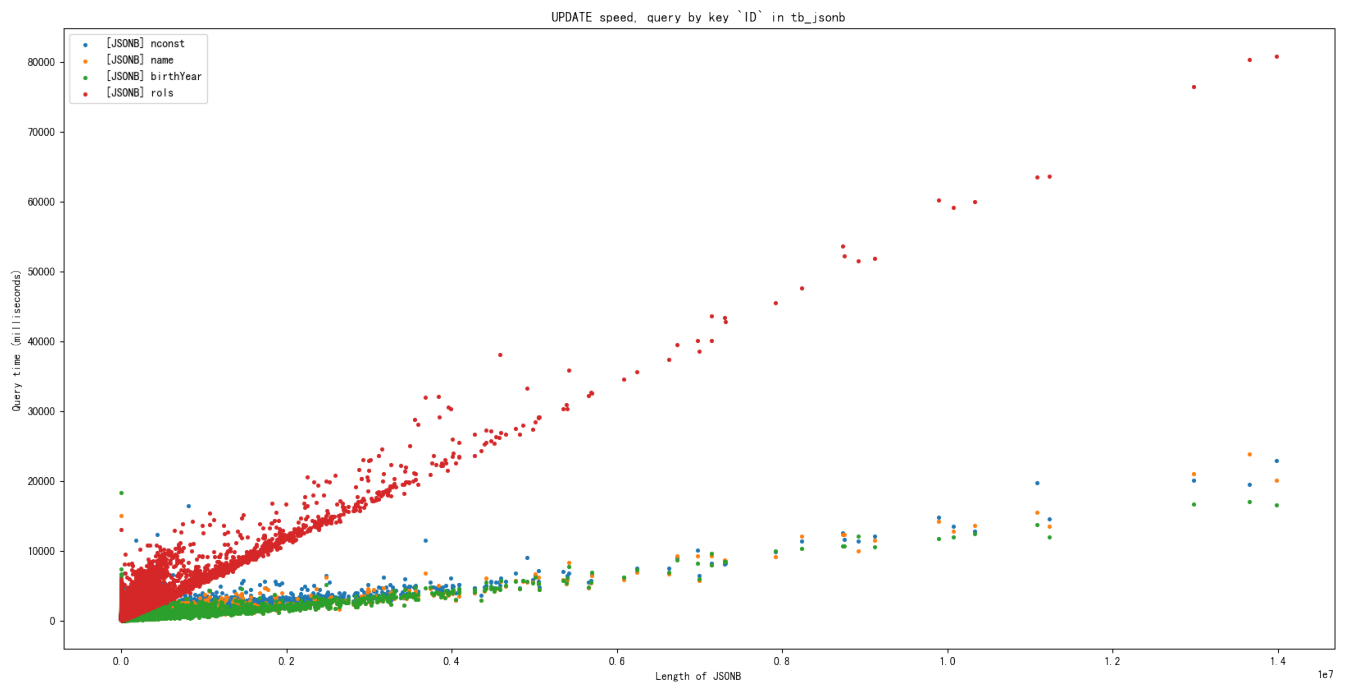
Сравнение скорости запросов к таблицам JSON и JSONB для всех полей:



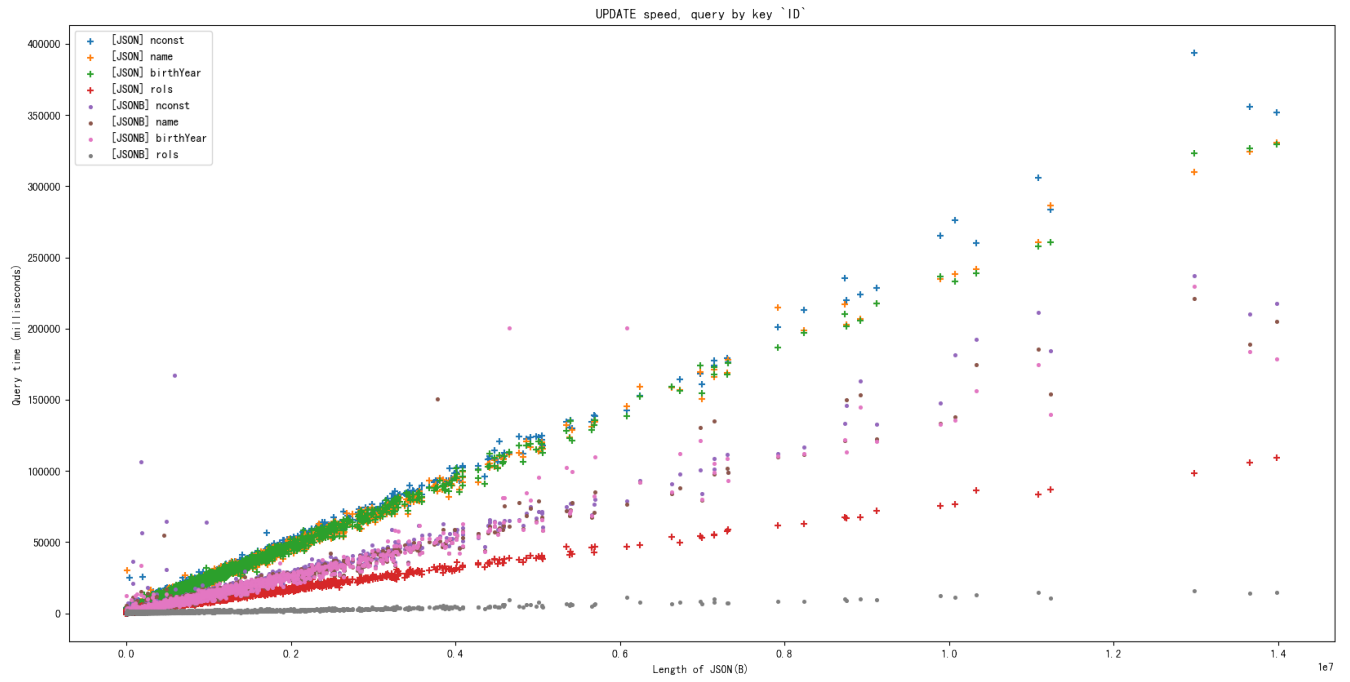
Сравнение времени UPDATE для каждого поля JSON в таблице tb_json :



Сравнение времени UPDATE для каждого поля JSONB в таблице tb_jsonb :



Сравнение скорости UPDATE к таблицам JSON и JSONB для всех полей :



1.2.2.2 TOAST

В PG страница является основной единицей хранения данных в файле, ее размер фиксирован и может быть задан только во время компиляции и не может быть изменен впоследствии, **размер по умолчанию составляет 8 КБ**.

Кроме того, **PG не позволяет хранить строку данных на разных страницах**. Для очень длинных строк данных PG инициирует TOAST, который сжимает или нарезает большие поля на несколько физических строк и сохраняет их в другой системной таблице (таблица TOAST), этот тип хранения называется **внерядным хранением**.

Для каждого поля таблицы в PG существует четыре стратегии TOAST.

- **PLAIN** - позволяет избежать сжатия и хранения вне ряда. Разрешается выбирать только те типы данных, для хранения которых не требуется политика TOAST (например, типы int), в то время как для таких типов, как текст, требующих длины хранения, превышающей размер страницы, эта политика недопустима.
- **EXTENDED** - позволяет сжимать и хранить вне ряда. Как правило, сначала он сжимается, а если он все еще слишком большой, то сохраняется вне очереди. Это политика по умолчанию для большинства типов данных, которые могут быть TOASTed.
- **EXTERNAL** - позволяет хранить данные вне ряда, но без сжатия. Это значительно ускоряет операции подстроки для полей типа text и bytea. Такие поля, как строки, которые работают с частью данных, могут достичь более высокой производительности при использовании этой политики, поскольку нет необходимости считывать всю строку данных и затем распаковывать ее.
- **MAIN** - позволяет сжимать, но не хранить вне линии. На практике, однако, внепоточное хранение активируется в крайнем случае, когда других методов (например, сжатия) недостаточно для гарантированного хранения больших данных. Поэтому правильнее будет сказать, что хранение вне ряда вообще не должно использоваться.

Просмотр политики TOAST для таблицы tb_jsonb:

Таблица TOAST имеет три поля:

- **chunk_id** -- используется для указания OID конкретного значения TOAST, которое можно интерпретировать как все строки с одинаковым значением chunk_id, образующие строку данных в поле TOAST исходной таблицы (в данном случае, блога).
- **chunk_seq** -- используется для указания позиции данных ряда в общем массиве данных.
- **chunk_data** -- фактические данные чанка

```

1 | postgres=# \d+ tb_jsonb;
2 |
3 |           Table "public.tb_jsonb"
4 | Column | Type   | Collation | Nullable | Storage | Compression | Stats target | Description
5 |-----+-----+-----+-----+-----+-----+-----+-----+-----
6 | iddata | integer |           | not null | plain   |              |              |
7 | imdata | jsonb   |           |          | extended |              |              |
8 |
9 | Indexes:

```

```

10  "ix_jsonb_iddata" btree (iddata)
11  Access method: heap
12
13
14
15  db_imdb=# select relname, relfilenode, reltoastrelid from pg_class where relname='tb_jsonb';
16  relname | relfilenode | reltoastrelid
17  -----+-----+-----
18  tb_jsonb |    162078   |    162082
19  (1 row)
20
21
22
23  db_imdb=# \d+ pg_toast.pg_toast_162078;
24  TOAST table "pg_toast.pg_toast_162078"
25    Column | Type | Storage
26  -----+-----+-----
27  chunk_id | oid  | plain
28  chunk_seq | integer | plain
29  chunk_data | bytea | plain
30  Owning table: "public.tb_jsonb"
31  Indexes:
32    "pg_toast_162078_index" PRIMARY KEY, btree (chunk_id, chunk_seq)
33  Access method: heap
34
35
36  postgres=# select * from pg_toast.pg_toast_162078;

```

Тестирование политики TOAST в Postgresql

1. [EXTENDED] Сравнить изменение объема БД для актера с малым кол-вом ролей

```

1 BEGIN;
2 SELECT * FROM tb_jsonb WHERE iddata=51989;
3 SELECT pg_table_size('tb_jsonb'); -- 1145241600 Byte
4 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1092 MB
5 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>' {name}' FROM tb_jsonb WHERE
iddata=51989; -- 202 Byte
6
7 UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{name}', "'Bf
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA'::jsonb) WHERE iddata=51989;
8
9 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>' {name}' FROM tb_jsonb WHERE
iddata=51989; -- 230 Byte
10 SELECT pg_table_size('tb_jsonb'); -- 1145241600 Byte
11 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1092 MB
12 ROLLBACK;

```

2. [EXTERNAL] Сравнить изменение объема БД для актера с малым кол-вом ролей

Политика TOAST изменена на EXTERNAL для отключения сжатия

```

1 postgres=# BEGIN;
2 postgres=# ALTER TABLE tb_jsonb ALTER imdata SET STORAGE EXTERNAL;
3
4 db_imdb=# \d+ tb_jsonb;
5
6          Table "public.tb_jsonb"
7  Column | Type   | Collation | Nullable | Storage | Compression | Stats target | Description
8  -----+-----+-----+-----+-----+-----+-----+-----
9  iddata | integer |           | not null | plain   |              |              |
10  imdata | jsonb   |           |          | external |              |              |
11
12 Indexes:
13   "ix_jsonb_iddata" btree (iddata)
14
15 Access method: heap

```

```

1 BEGIN;
2 ALTER TABLE tb_jsonb ALTER imdata SET STORAGE EXTERNAL;
3
4 SELECT * FROM tb_jsonb WHERE iddata=51989;
5 SELECT pg_table_size('tb_jsonb'); -- 1202307072 Byte
6 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1147 MB
7 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>' {name}' FROM tb_jsonb WHERE
iddata=51989; -- 202 Byte
8 select count(*) from pg_toast.pg_toast_162078; -- 511805
9

```

```

10 UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{name}', "'Bf
    AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA'"::jsonb) WHERE iddata=51989;
11
12 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>' {name}' FROM tb_jsonb WHERE
    iddata=51989; -- 230 Byte
13 SELECT pg_table_size('tb_jsonb'); -- 1202307072 Byte
14 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1147 MB
15 select count(*) from pg_toast.pg_toast_162078; -- 511805
16 ROLLBACK;

```

3. [EXTENDED] Сравнить изменение объема БД для актера с большим кол-вом ролей

```

1 BEGIN;
2 SELECT pg_table_size('tb_jsonb'); -- 1145241600 Byte
3 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1092 MB
4 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>' {name}' FROM tb_jsonb WHERE
    iddata=3789; -- 4034997 Byte
5
6 UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{name}', "'David
    AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA'"::jsonb) WHERE iddata=3789;
7
8 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>' {name}' FROM tb_jsonb WHERE
    iddata=3789; -- 4035007 [+10] Byte
9 SELECT pg_table_size('tb_jsonb'); -- 1149378560 [+4136960] Byte
10 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1096 MB [+3.9453125 MB]
11 ROLLBACK;

```

4. [EXTERNAL] Сравнить изменение объема БД для актера с большим кол-вом ролей

Политика TOAST изменена на EXTERNAL для отключения сжатия

```

1 postgres=# BEGIN;
2 postgres=# ALTER TABLE tb_jsonb ALTER imdata SET STORAGE EXTERNAL;
3
4 db_imdb=# \d+ tb_jsonb;
5
6          Table "public.tb_jsonb"
7  Column | Type   | Collation | Nullable | Storage | Compression | Stats target | Description
8  -----+-----+-----+-----+-----+-----+-----+-----
9  iddata | integer |           | not null | plain   |              |              |
10  imdata | jsonb   |           |          | external |              |              |
11
12 Indexes:
13     "ix_jsonb_iddata" btree (iddata)
14
15 Access method: heap

```

Мы можем видеть, что объем данных значительно больше без включенного сжатия.


```

1 BEGIN;
2 ALTER TABLE tb_jsonb ALTER imdata SET STORAGE EXTERNAL;
3
4 SELECT pg_table_size('tb_jsonb'); -- 1186045952 Byte
5 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1131 MB
6 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>'{'name}' FROM tb_jsonb WHERE
iddata=3789; -- 4034997 Byte
7
8 UPDATE tb_jsonb SET imdata=jsonb_set(imdata::jsonb, '{name}', '"David
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAA":jsonb)' WHERE iddata=3789;
9
10 SELECT iddata, pg_column_size(imdata) , imdata, imdata->>'{'name}' FROM tb_jsonb WHERE
iddata=3789; -- 15845197 [+11810200] Byte
11 SELECT pg_table_size('tb_jsonb'); -- 1202307072 [+16261120] Byte
12 SELECT pg_size_pretty(pg_table_size('tb_jsonb')); -- 1147 MB [+16 MB]
13 ROLLBACK;

```

1.2.2.2.1 Результат

- Если политика разрешает сжатие, TOAST предпочитает сжатие.
- Хранение вне ряда включается, когда объем данных превышает примерно 2 КБ, независимо от того, сжаты они или нет.
- Изменение политики TOAST **не** повлияет на способ хранения существующих данных.

2. Приложение

2.1 Адрес репозитория GitHub

Все исходные данные, выводы и код можно найти на GitHub, адрес репозитория:

<https://github.com/NekoSilverFox/PostgreSQL-SPbSTU>