



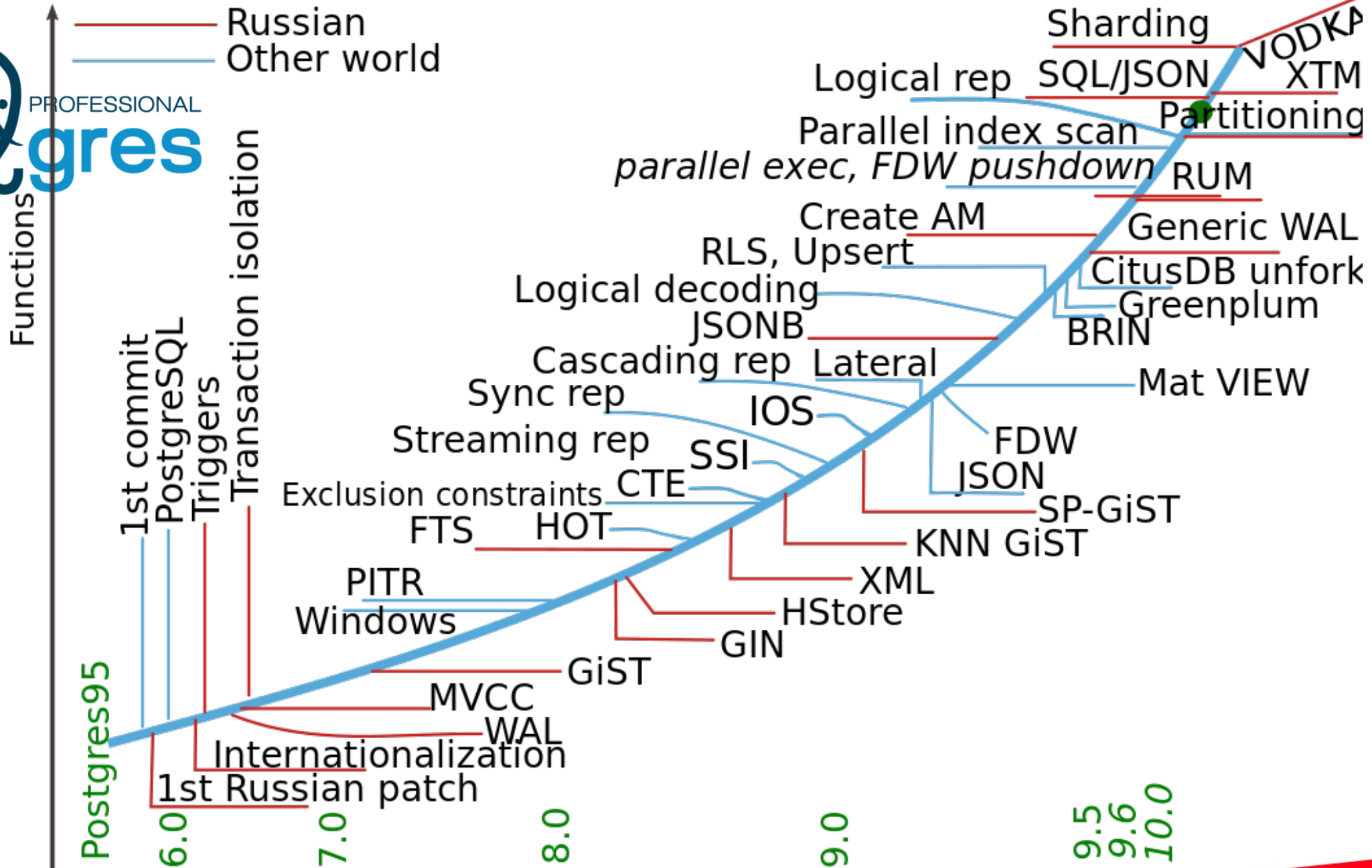
# Оптимизация high-contention write в PostgreSQL

Олег Бартунов, Александр Коротков, Иван Панченко



HighLoad<sup>++</sup>

Профессиональная конференция  
разработчиков высоконагруженных  
систем





How Postgres is good  
as NoSQL database ?

Возможности богатые, а так  
кто его знает ...

SQL/JSON 2016 ! Пора наконец  
померять производительность!

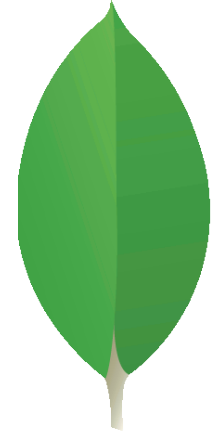
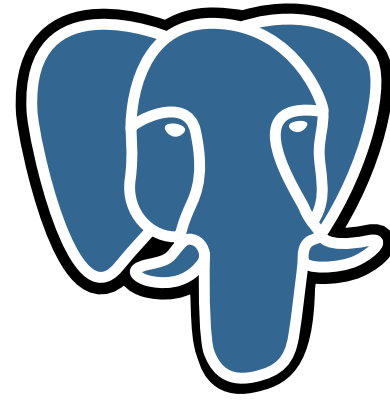


# Benchmarking NoSQL Postgres

- Existing benchmarks were homemade by postgres people
- People tend to believe independent and «scientific» benchmarks
  - Reproducible
  - More databases
  - Many workloads
  - Open source



# YCSB Benchmark



- Yahoo! Cloud Serving Benchmark -  
<https://github.com/brianfrankcooper/YCSB/wiki>
- De-facto standard benchmark for NoSQL databases
- Scientific paper «Benchmarking Cloud Serving Systems with YCSB»  
<https://www.cs.duke.edu/courses/fall13/cps296.4/838-CloudPapers/ycsb.pdf>
- We run YCBS for Postgres master, MongoDB 3.4.2
  - 1 server with 72 cores, 3 TB RAM, 2 TB SSD for clients
  - 1 server with 72 cores, 3 TB RAM, 2 TB SSD for database
  - 10Gbps switch

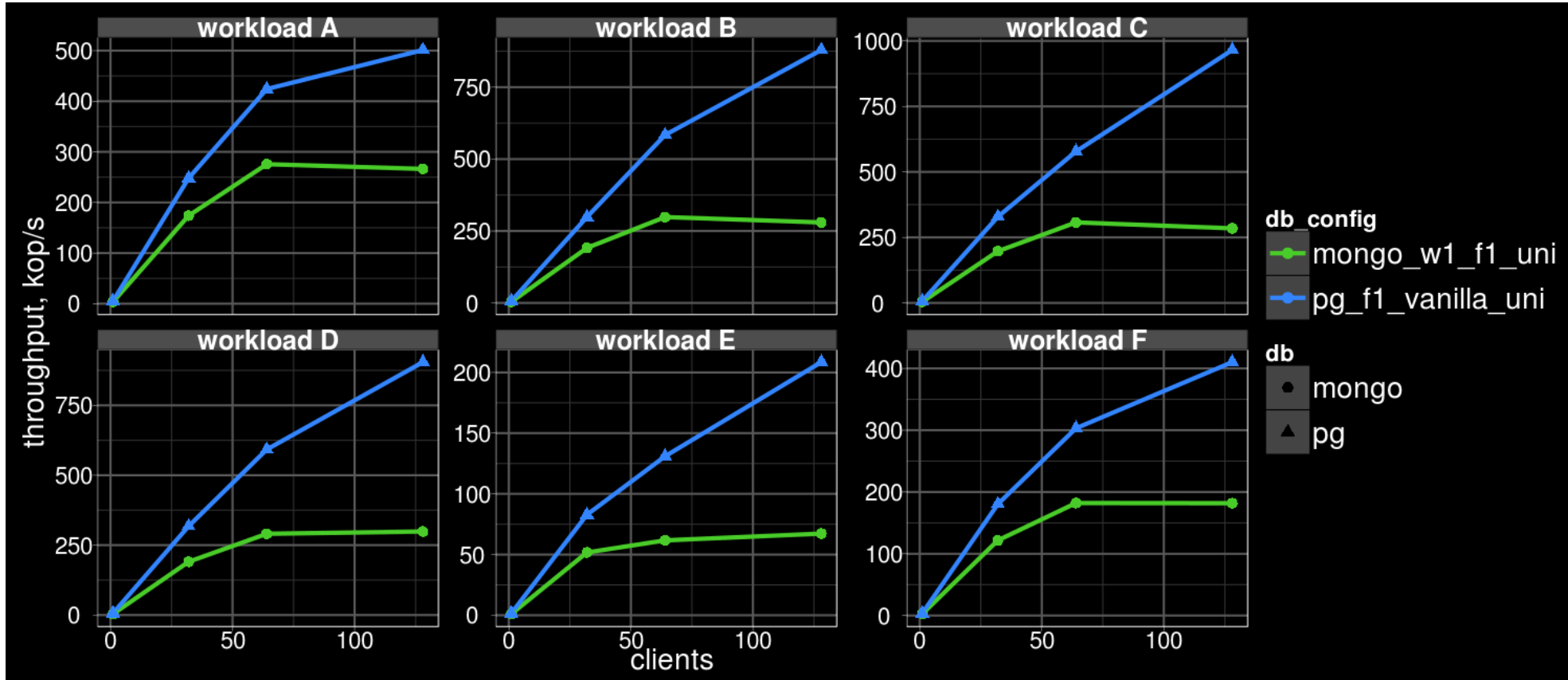
# YCSB Benchmark: Core workloads

- Workload A: Update heavy - a mix of 50/50 reads and writes
- Workload B: Read mostly - a 95/5 reads/write mix
- Workload C: Read only — 100% read
- Workload D: Read latest - new records are inserted, and the most recently inserted records are the most popular
- Workload E: Short ranges - short ranges of records are queried
- Workload F: Read-modify-write - the client will read a record, modify it, and write back the changes
- All (except D) workloads uses Zipfian distribution for record selections

# YCSB Benchmark: details (1)

- #1: Postgres 10, synchronous commit=off  
Mongodb 3.4.2 (w1, j0) — 1 mln. Rows, #clients <= 128
- #2: Postgres 10, synchronous commit=off  
Mongodb 3.4.2 (w1, j0) — 1 mln. Rows, #clients <= 1000
- #3: Postgres 10, synchronous commit=on  
Mongodb 3.4.5 (w1, j1)
- We tested:
  - Functional btree index for jsonb, jsonbc, sqljson
  - Mongodb (wiredtiger with snappy compression)
  - Return a whole json, just one field, small range

# Uniform distribution of queries: #1





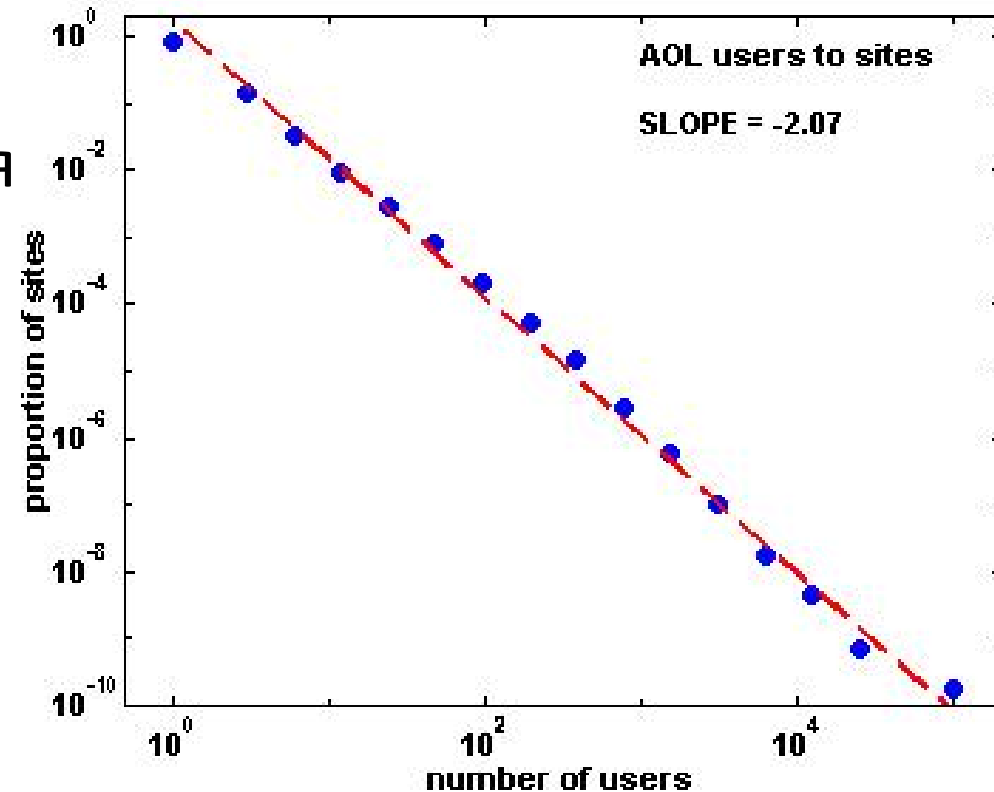
# Распределение Зипфа и PostgreSQL

- Разработчики в основном пользуются pgbench. PostgreSQL – pgbench optimized DBMS.
- YCSB предлагает использовать распределение Зипфа
- В pgbench не было распределения Зипфа, исправляемся

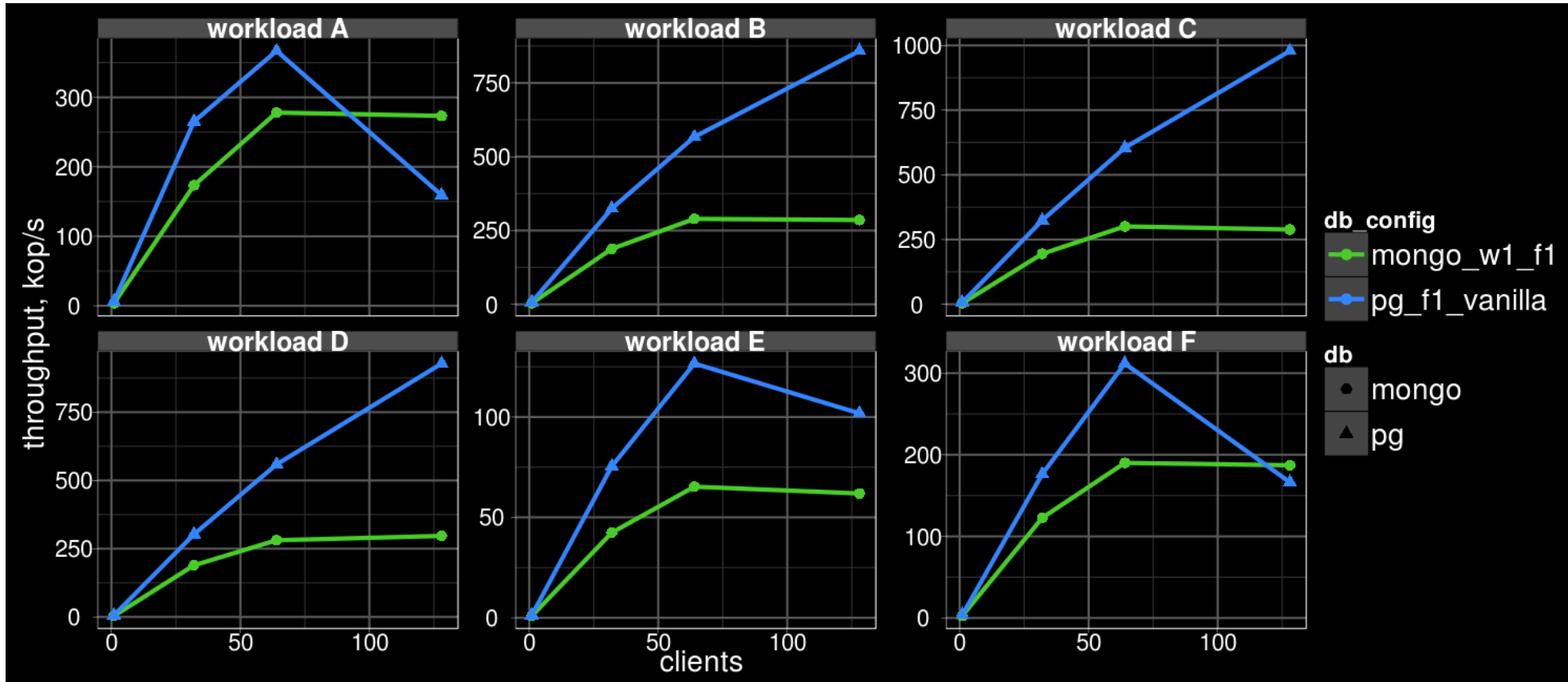
<https://www.postgresql.org/message-id/flat/BF3B6F54-68C3-417A-BFAB-FB4D66F2B410@postgrespro.ru>

# Что это за распределение Зипфа

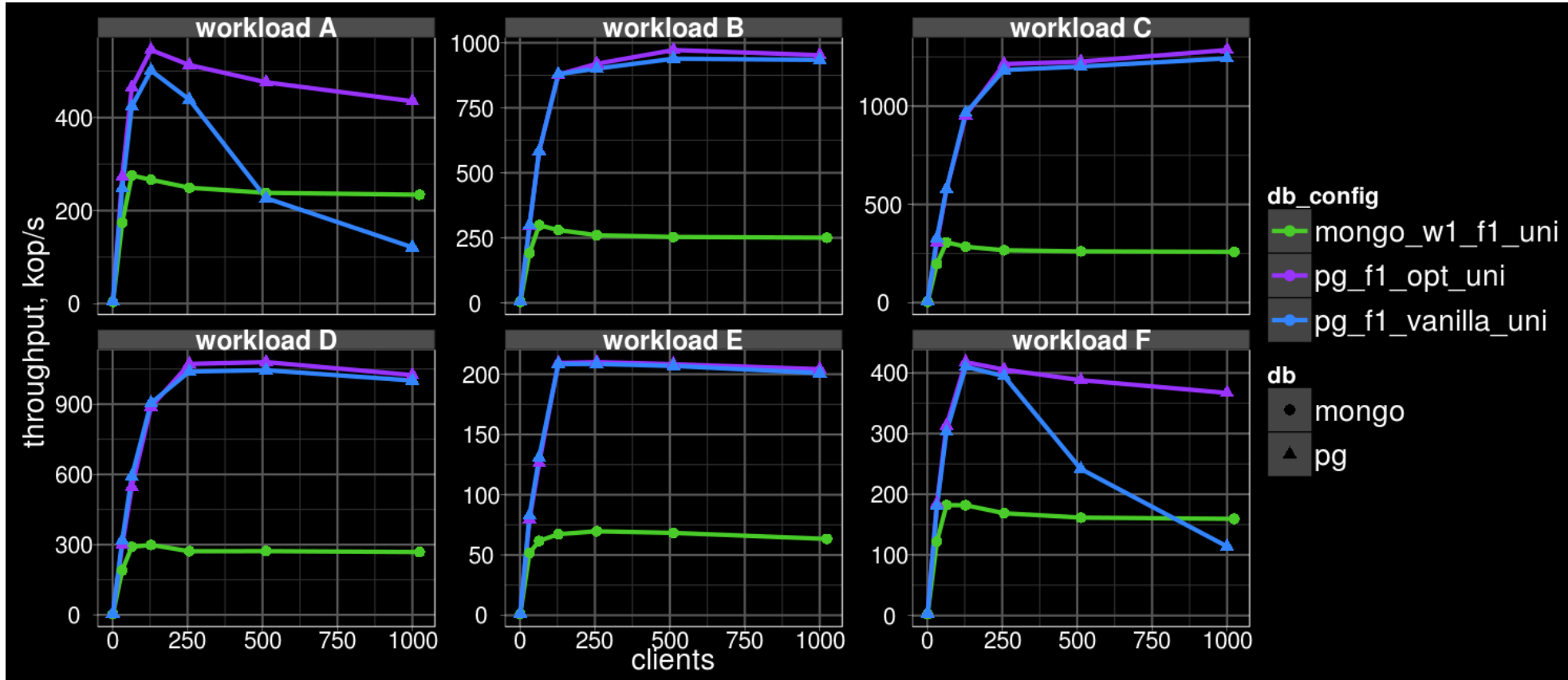
- Частота обратно пропорциональна рангу значения
- Т.е. второе по частоте значение встречается в 2 раза реже чем 1е
- Третье — в три раза реже, и т. д.
- Эмпирически выведено из анализа частотности слов
- Хорошо аппроксимирует другие жизненные ситуации



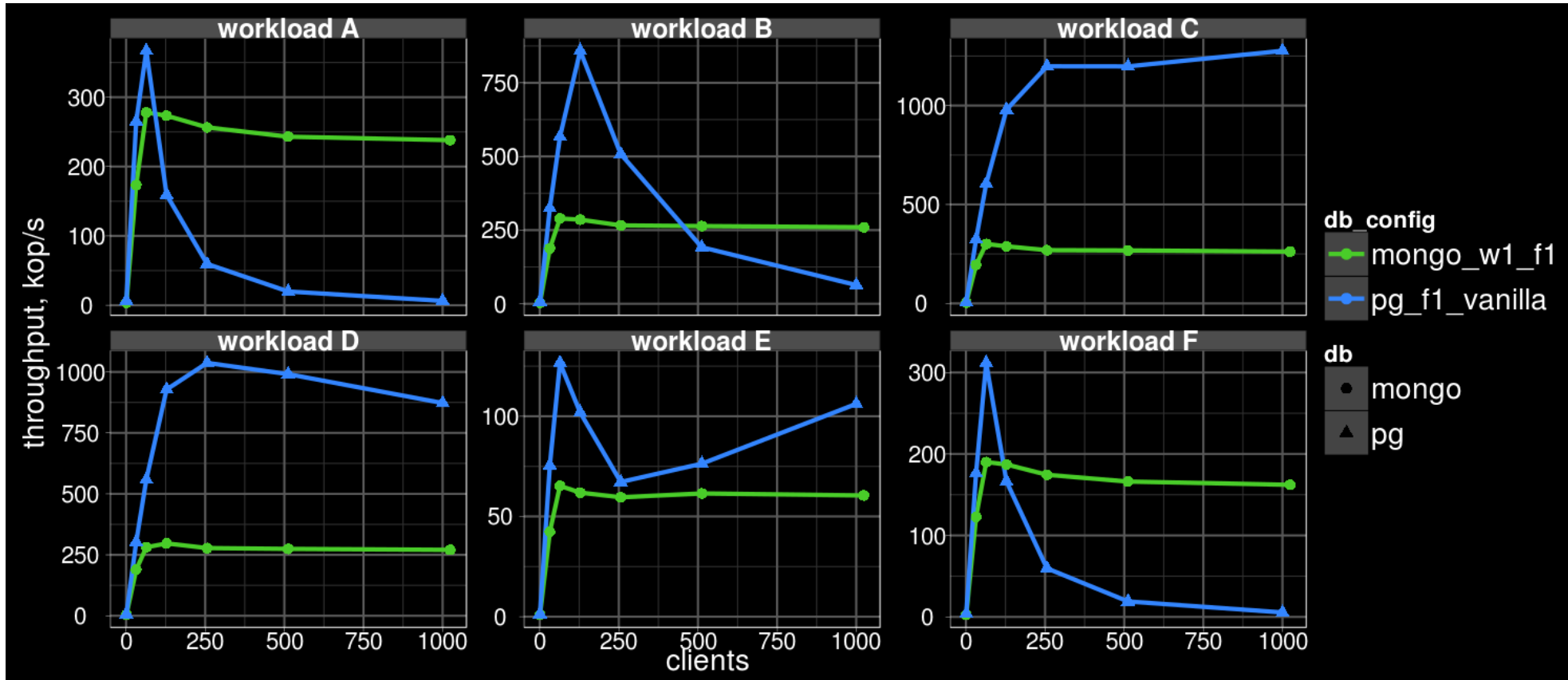
# Zipfian distributions of queries #1



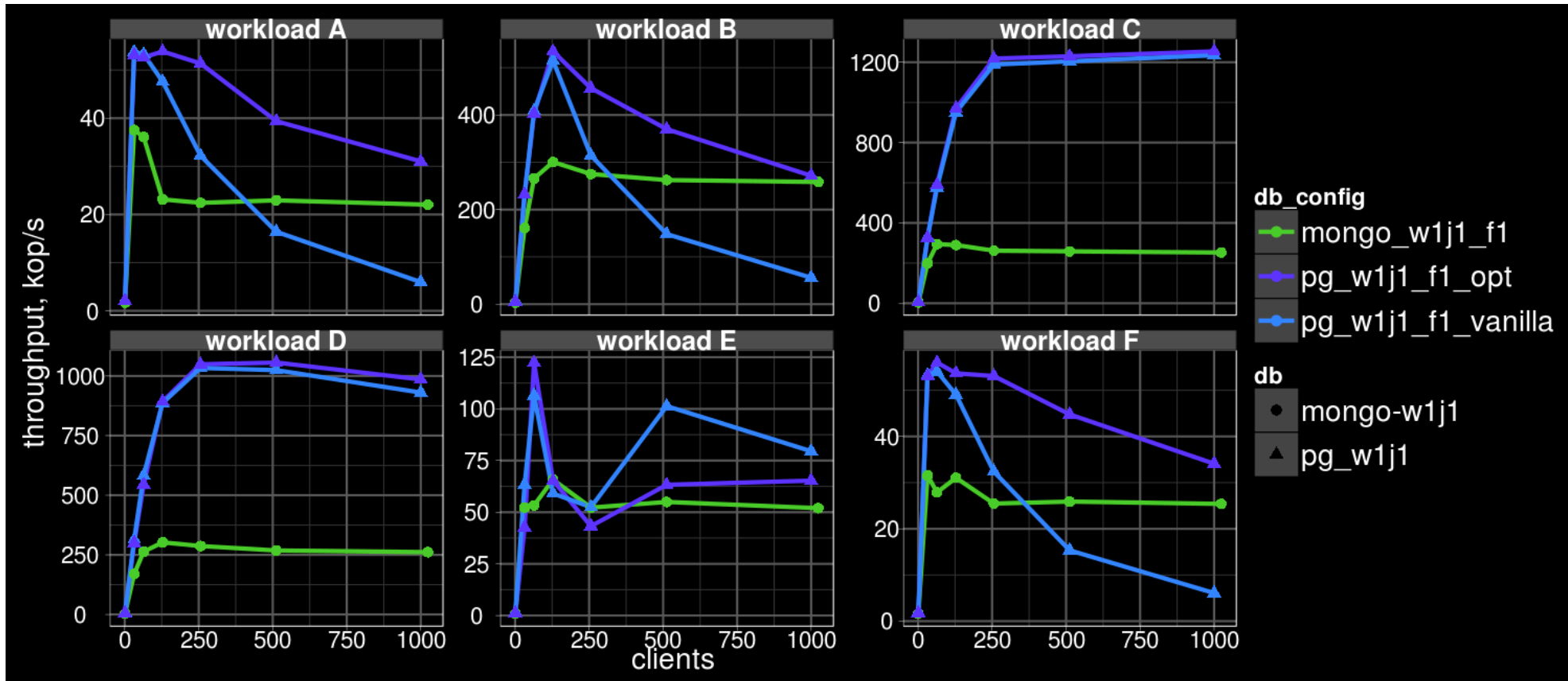
# Uniform distribution of queries: #2



# Zipfian distributions of queries: #2



# Zipfian distributions of queries: #3



## Выводы

- Сочетание «`synchronous_commit = off`», zipfian distribution и большой конкурентности ( $> 100$ ) не очень хорошо для постгреса (обновление)
- Даже для однородного распределения при большом количестве клиентов ( $> 500$ ) производительность обновления постгреса деградирует
- Персистентность «`synchronous_commit = on`» немного помогает постгресу конкурировать с MongoDB (j1), но не кардинально



Ого, а постгрес  
ничего !

Для «вебовской» нагрузки постгрес  
сливает при большой конкуретности

Надо разбираться !

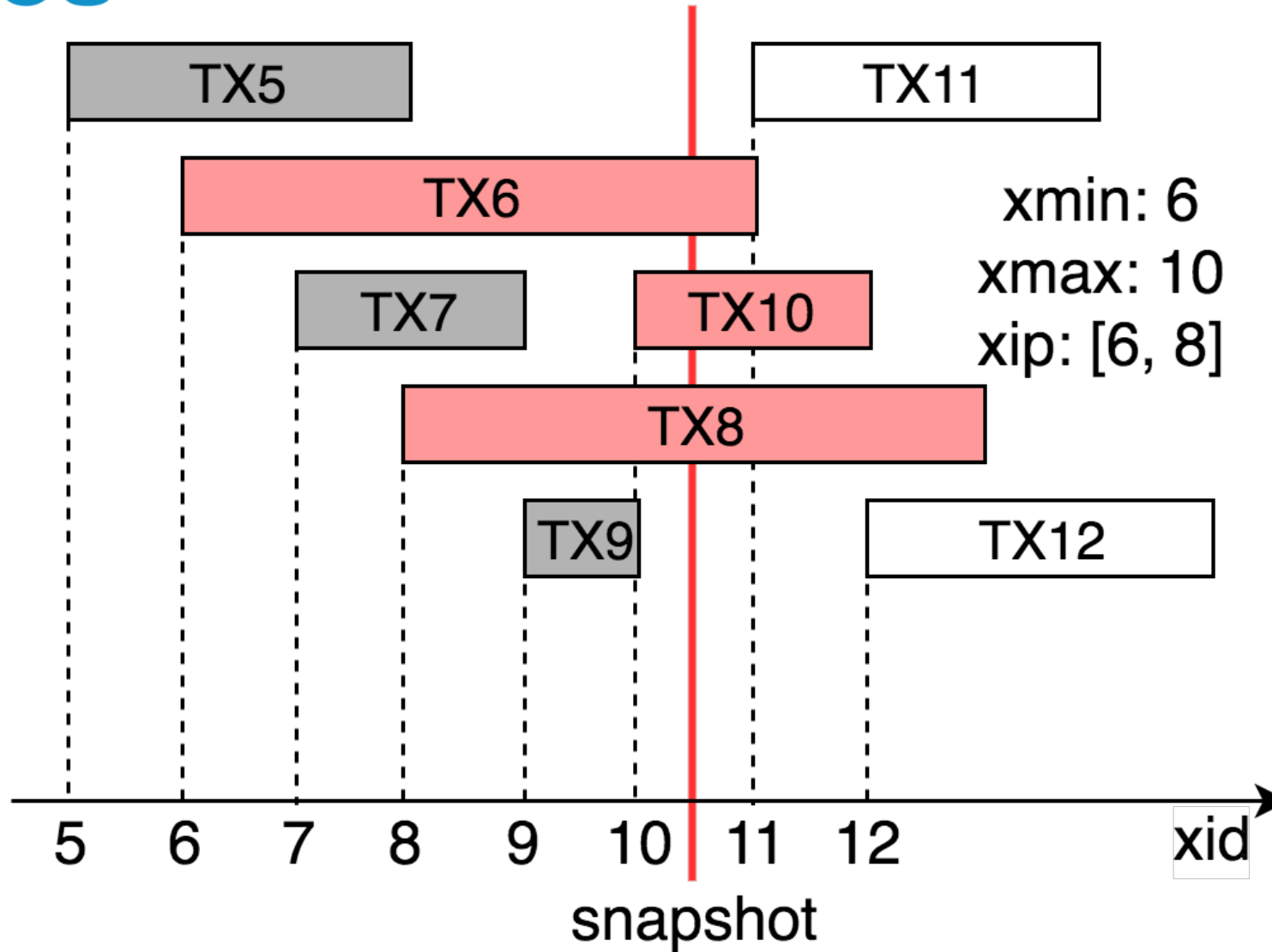


Isolation – на выполнение данной транзакции не должны оказывать влияние другие транзакции, которые работают параллельно с ней.

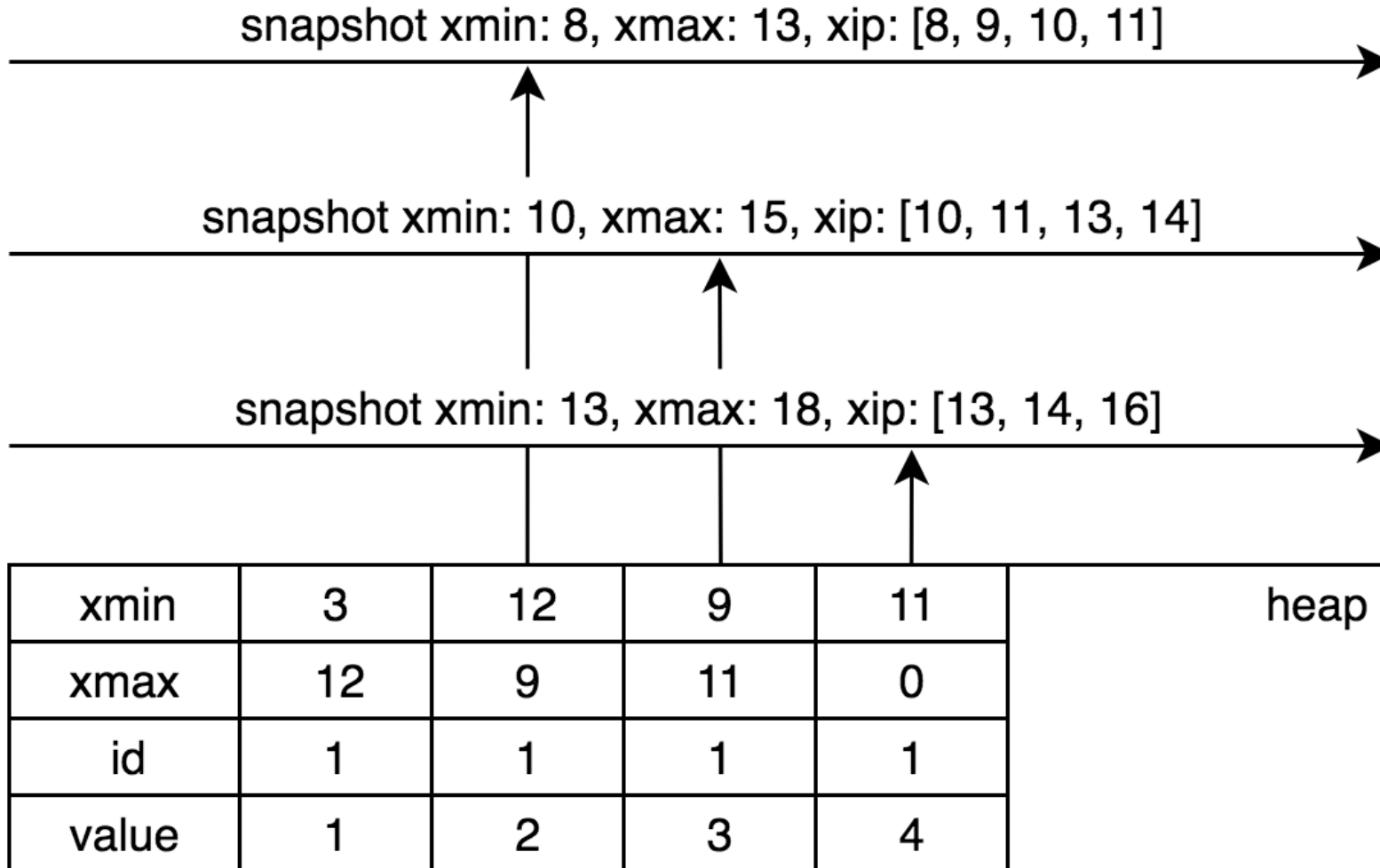
Для того, чтобы выполнить свойство isolation, но при этом читающие и пишущие транзакции не блокировали друг друга, был придуман MVCC. При MVCC старые версии строк сохраняются для того, чтобы читающие транзакции могли видеть их.

Snapshot – это состояние базы данных между двумя транзакциями. Для того, чтобы идентифицировать снапшот, нужно уметь отличать, какие транзакции в прошлом, а каким – в будущем. В PostgreSQL снапшот задаётся с помощью xmin, xmax и массива xip.

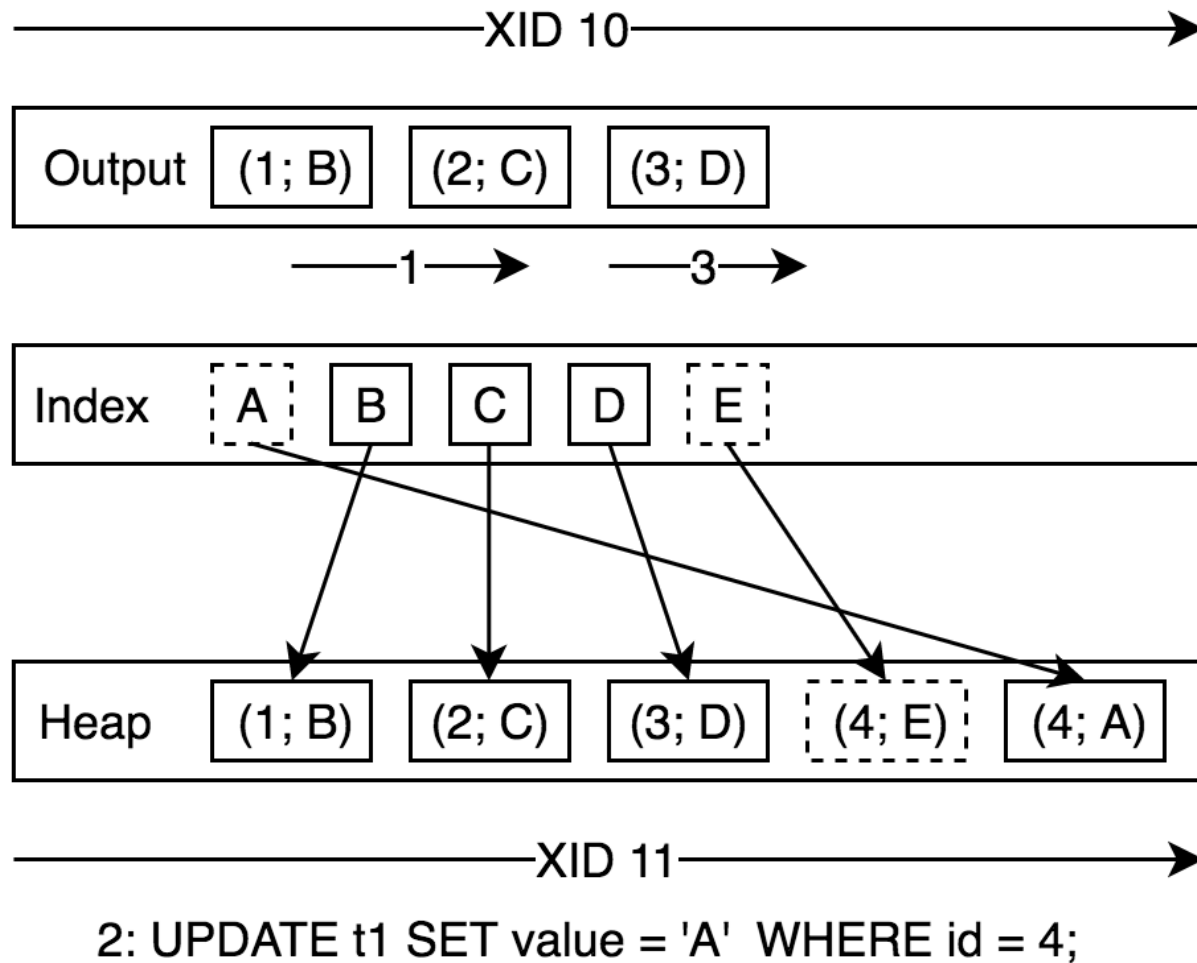
# Что такое snapshot



# Использование snapshot'ов



# Read committed vs snapshots

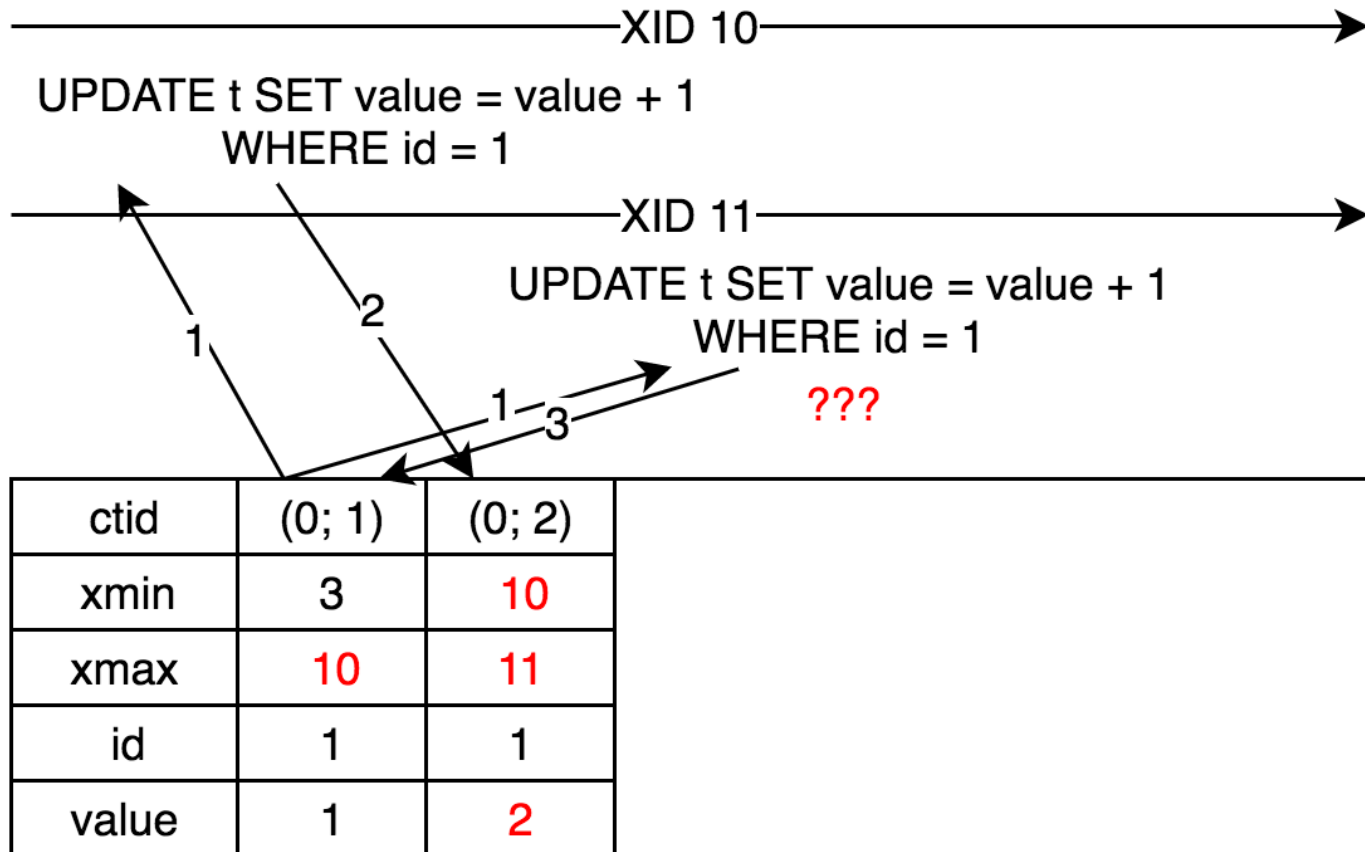


- 1) XID 10 читает туплы (1; B) и (2; C).
- 2) XID 11 обновляет тупл (4; E) на (4; A) и коммитится.
- 3) XID 10 читает тупл (3; D) , а (4; E) ему уже не видим, т. к. он уже был обновлён.

В итоге мы потеряли строку. Чтобы так не было нужно вначале сделать

- 0) XID 10 берёт снапшот для читающего запроса.

# Contended update



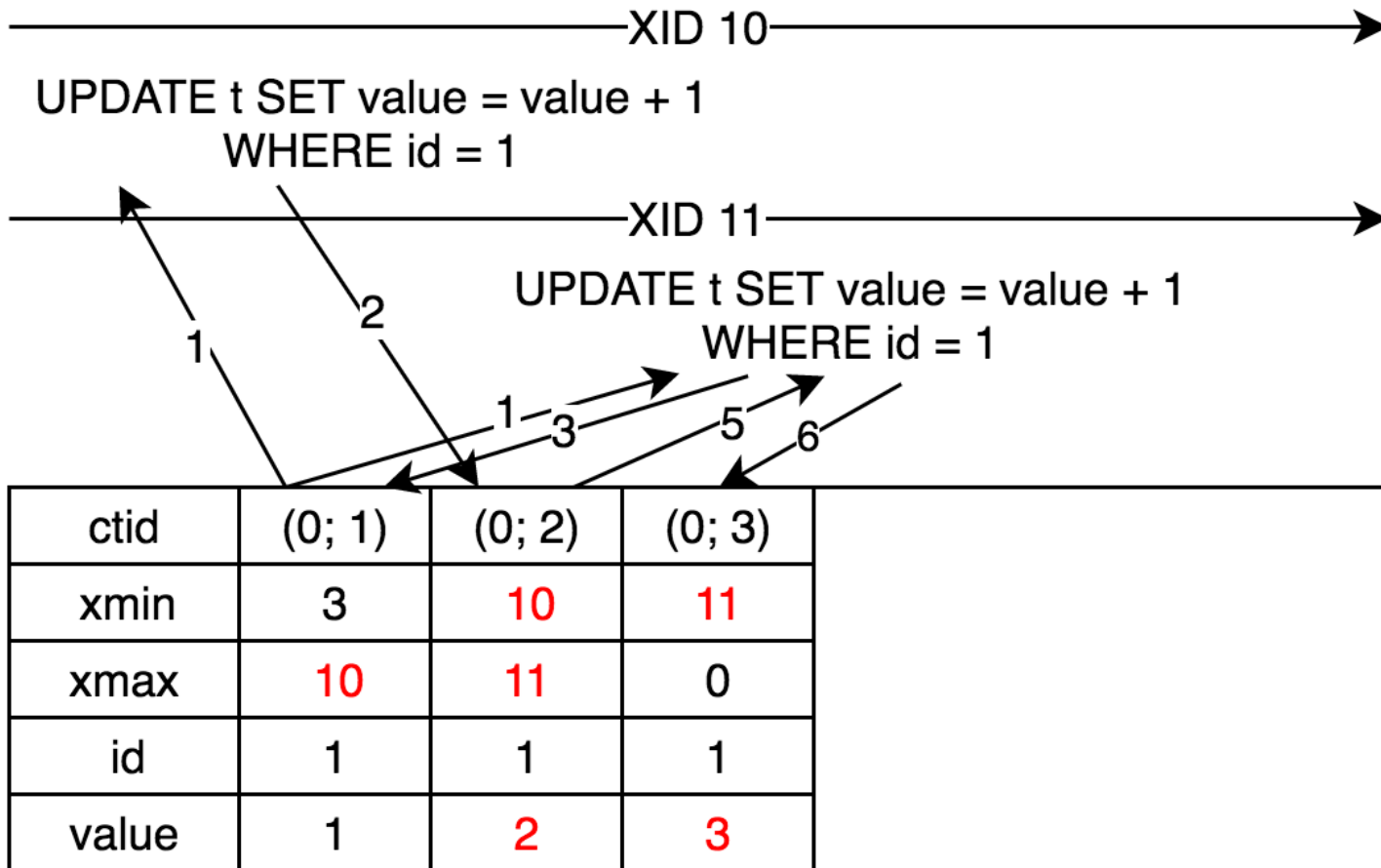
- 1) Обе транзакции прочитали в соответствии со своими снапшотами одну и ту же версию строки ctid = (0;1). И подготовили обновлённые версии с value = 2.
- 2) xid 10 первая обновляет строку и добавляет версию ctid = (0;2) с value = 2.
- 3) xid 11 тоже пытается обновить строку ctid = (0;1), но видит что она уже обновлена... Что делать дальше?

isolation level  $\geq$  repeatable read

```
ERROR: could not serialize access due to concurrent update
```

Если нам приходится обновлять строку, которую уже успела обновить другая параллельная транзакция, то это значит что её результат отразится на результатах нашей транзакции. Таким образом, для строгого выполнения свойства `isolation` нам ничего не остаётся как выдать ошибку (и попробовать нашу транзакцию ещё раз).

# Contended update: read committed



- 1) Обе транзакции прочитали в соответствии со своими снапшотами одну и ту же версию строки ctid = (0;1). И подготовили обновлённые версии с value = 2.
- 2) xid 10 первая обновляет строку и добавляет версию ctid = (0;2) с value = 2.
- 3) xid 11 тоже пытается обновить строку ctid = (0;1), но видит что она уже обновлена и ждёт завершения xid 10.
- 4) xid 10 коммитится.
- 5) xid 11 просыпается, находит последнюю версию строки ctid = (0;2) и лочит её.
- 6) xid 11 перевычисляет обновлённую версию строки (value = 3) и вставляет её в heap.

## Как мы лочили tuple в PostgreSQL 9.4

- 1) Дожидаемся, пока транзакция, записанная в хтах, закончится.
- 2) Ищем последний tuple в цепочке update'ов.
- 3) Пытаемся записаться свой xid в tuple хтах. Если кто-то уже сделал это до нас, то переходим к шагу 1.

При большом contention'е всё было совсем плохо.

Поправили в [0e5680f4](#).

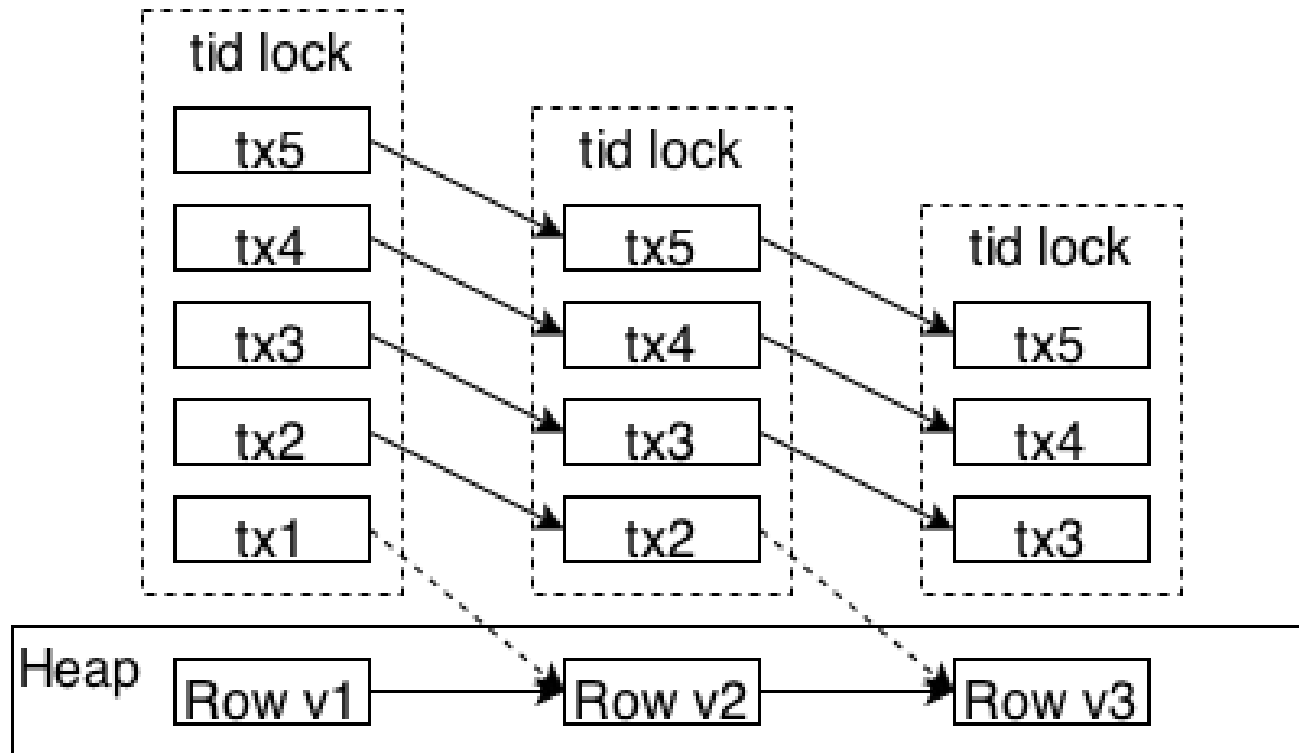


## Как мы лочим tuple в PostgreSQL 9.5+

- 1) Берем Lock на tuple id.
- 2) Дожидаемся, пока транзакция, записанная в xmax, закончится.
- 3) Ищем последний tuple в цепочке update'ов.
- 4) Записываем свой xid в tuple xmax.

Стало лучше – транзакции стали выстраиваться в очередь!

## Стало лучше, но...



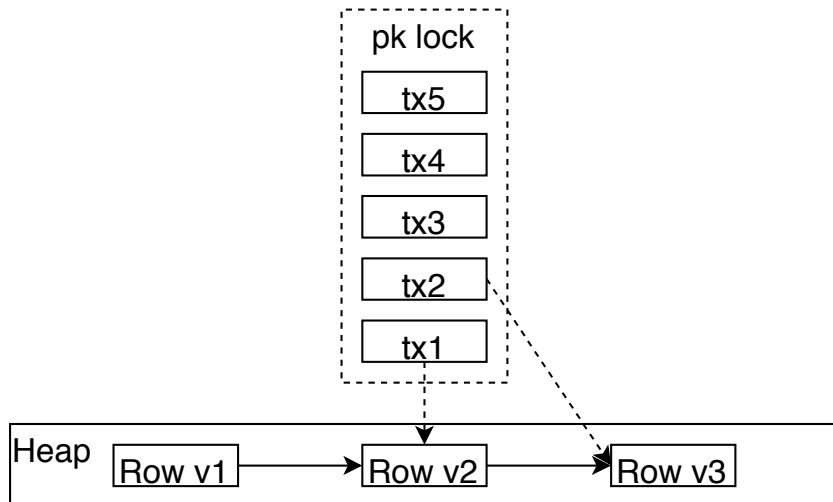
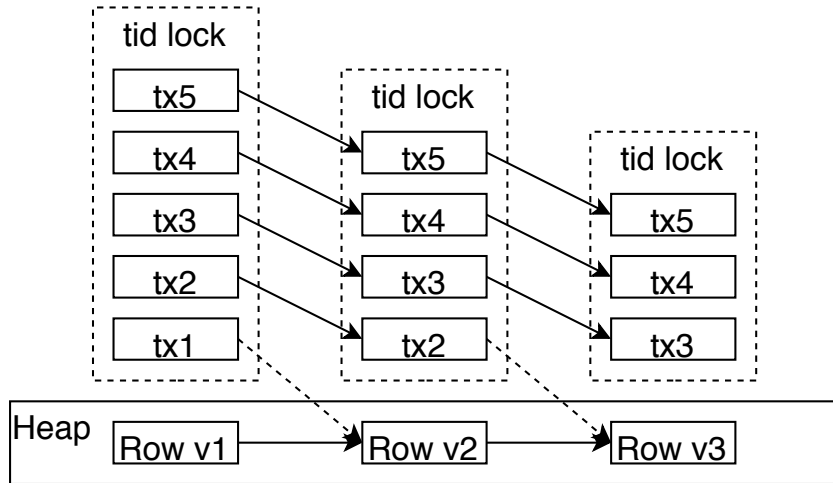
- tuple id меняется и всё очередь постоянно переезжает.
- Получается  $O(n^2)$  операция взятия/отпускаения лока.
- Heavy weight lock сам по себе не дешёвый. Включается в себя заглядывание в хэш-таблицу в shared memory.

# Виды lock'ов

	spinlock	LWLock	HWLock
Заранее невозможно выделить shared memory под каждый lock	—	—	+
Число уровней блокировки	1	2	8
Очередь блокировок	—	+	+
Deadlock detection	—	—	+

Для tuple lock подходит только heavy-weight lock...

# Оптимизация: primary key lock



- В отличии от tuple id, значение primary key, как правило, не меняется.
- В итоге очередь на update одна, не переезжает.
- Нужна уметь получать хорошее (и достаточно длинное) хэшированное представление primary key, чтобы исключить коллизии.

# Оптимизации для high-contention write

- Primary key lock during update
- Низкоуровневые оптимизации

Small improvement to compactify\_tuples

<https://www.postgresql.org/message-id/flat/3c6ff1d3a2ff429ee80d0031e6c69cb7>

Fix performance degradation of contended LWLock on NUMA

<https://www.postgresql.org/message-id/flat/2968c0be065baab8865c4c95de3f435c@postgrespro.ru>

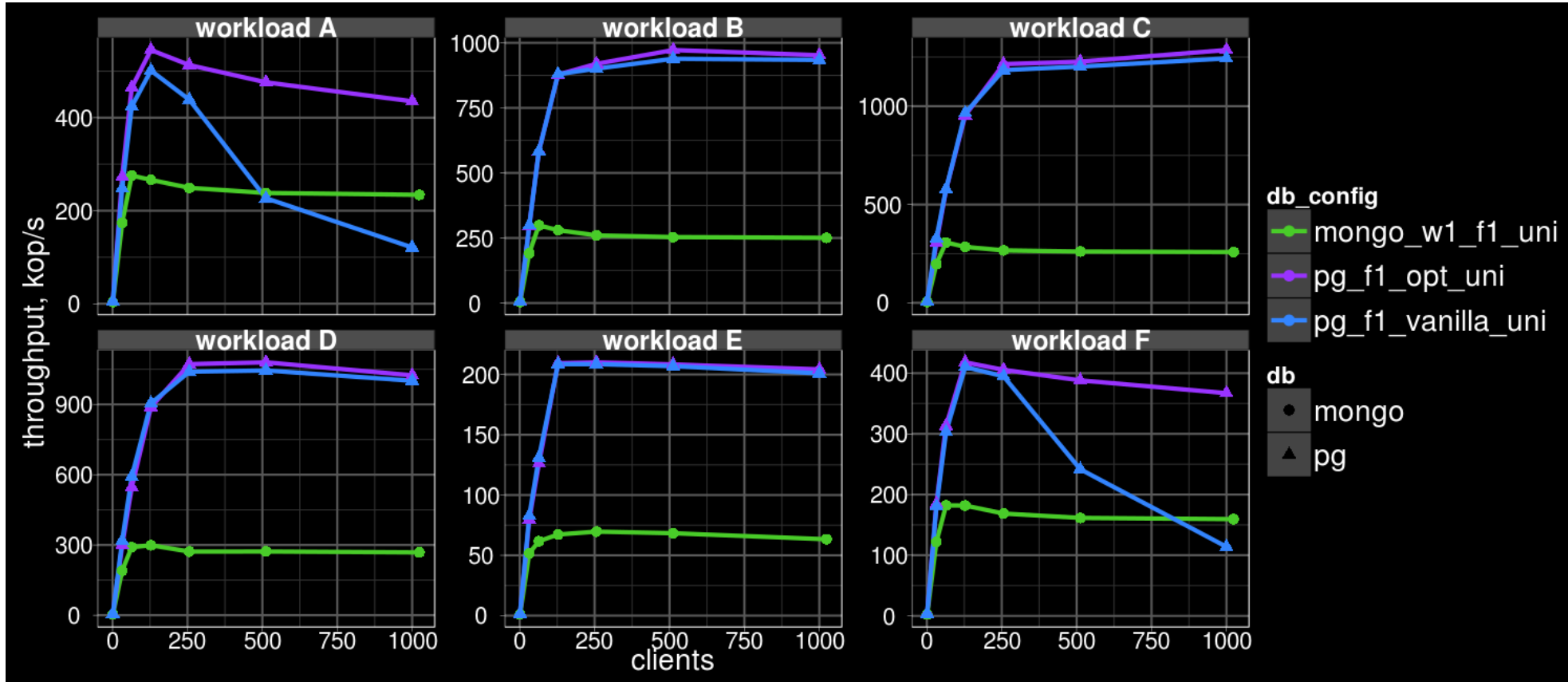


Что-то сложно все ...

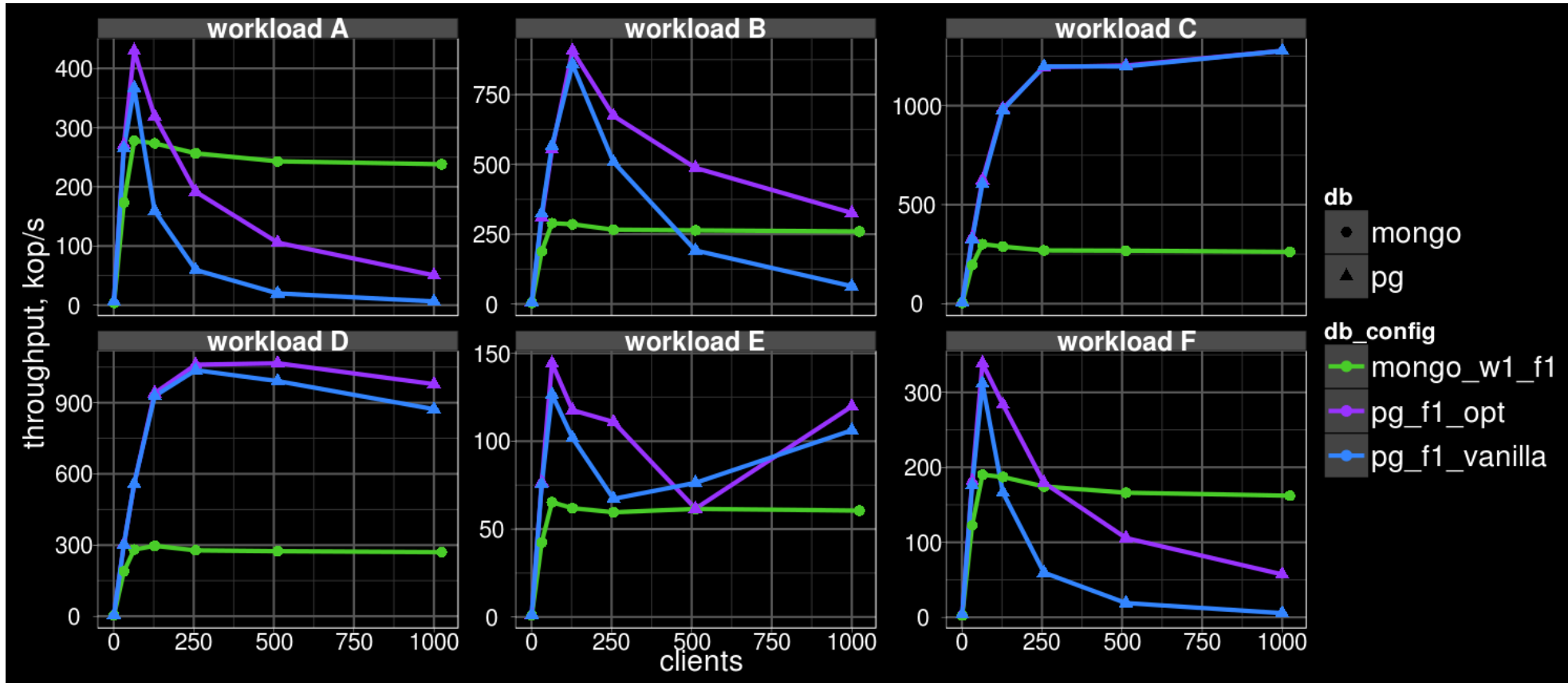
Это упрощенная картина :)

Давай результаты !

# Uni with postgres optimized — GOOD !

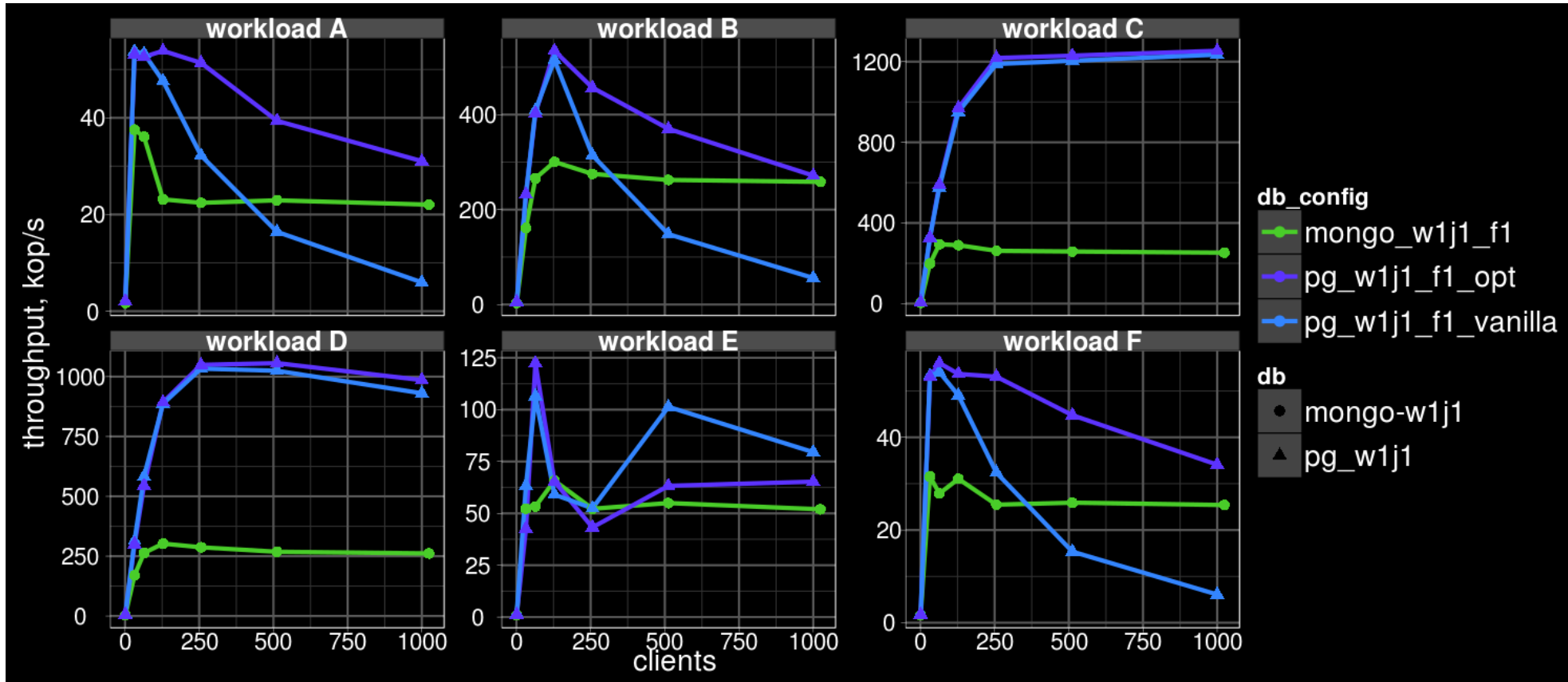


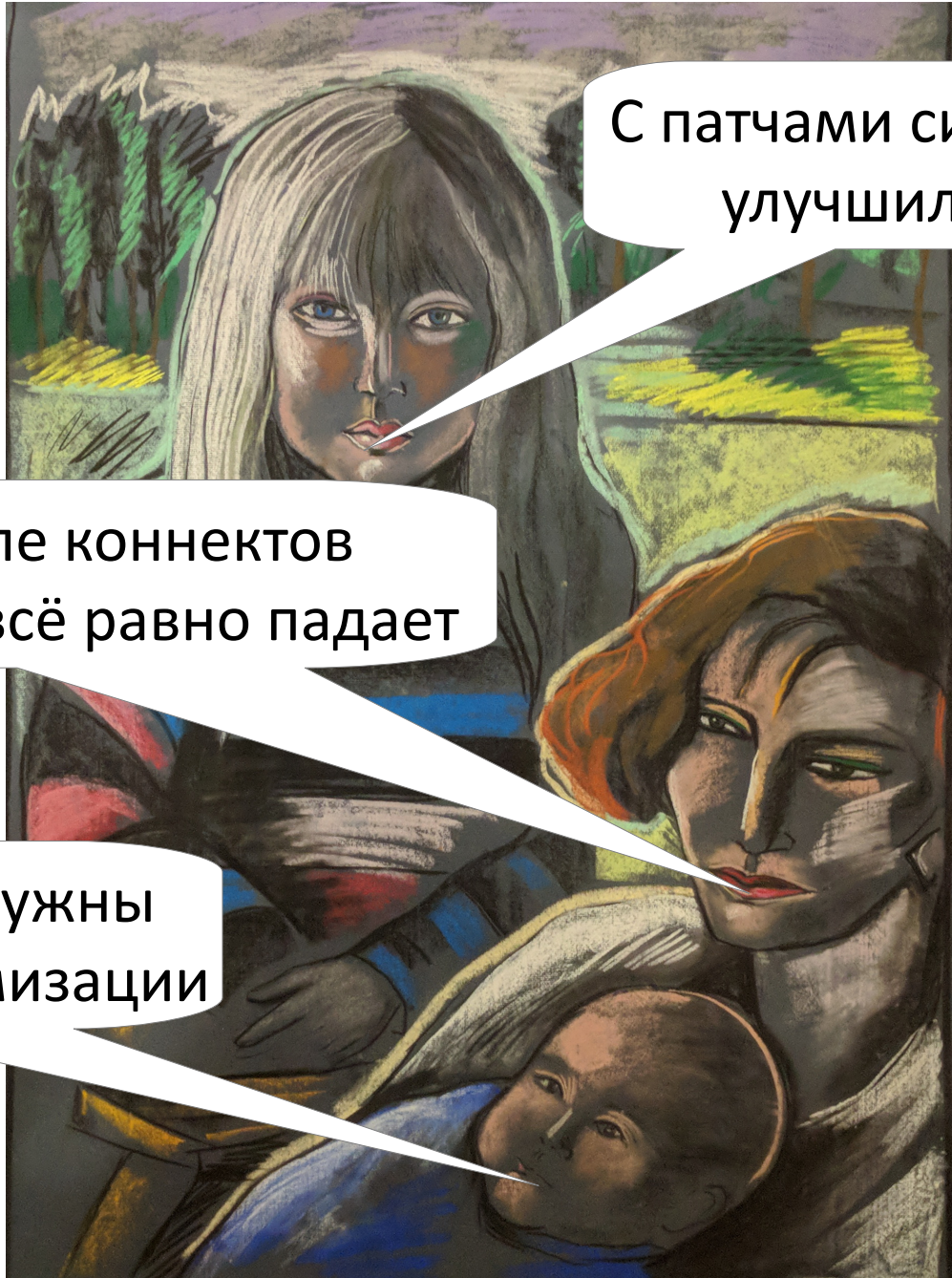
# Zipf with postgres optimized — So-So :(





# Persistent case (sync. Commit) — Good :)





С патчами ситуация  
улучшилась

При большом числе коннектов  
производительность всё равно падает

Проблема остаётся, нужны  
дополнительные оптимизации

## Что делать с high-contention write?

- Для большинства приложений текущей производительности PostgreSQL более чем достаточно.
- По-возможности избегать high-contention write.
- Ждать улучшений :)

# Как избежать high-contention write? (1/4)

```
CREATE TABLE post (  
    id bigserial primary key,  
    author text,  
    title text,  
    body text,  
    views_count bigint,  
    comments_count bigint  
);
```

```
CREATE TABLE post_comment (  
    id bigserial primary key,  
    post_id bigint  
        REFERENCES post (id),  
    author text,  
    body text  
);
```

## Как избежать high-contention write? (2/4)

`post.views_count`

- Строгая consistency, как правило, не нужна
- Накапливаем просмотры в приложении или in-memory БД
- Периодически по cron обновляем значения `post.views_count`

## Как избежать high-contention write? (3/4)

```
CREATE TABLE post (  
    id bigserial primary key,  
    author text,  
    title text,  
    body text,  
    views_count bigint,  
    comments_count bigint,  
    last_snapshot txid_snapshot  
);
```

```
CREATE TABLE post_comment (  
    id bigserial primary key,  
    post_id bigint  
        REFERENCES post (id),  
    author text,  
    body text,  
    insert_xid bigint
```

## Как избежать high-contention write? (4/4)

По trigger'у устанавливаем `post_comment.insert_xid = txid_current()`.  
По cron'у обновляем счётчик комментариев.

```
UPDATE post
SET comments_count = views_count + (
    SELECT count(*) FROM post_comment pc
    WHERE pc.post_id = post.id AND
          NOT txid_visible_in_snapshot(pc.insert_xid,
                                       post.last_snapshot)
),
last_snapshot = txid_current_snapshot();
```

## Выводы

- PostgreSQL показывает хорошие результаты на YCSB benchmark, напрямую конкурируя с NoSQL решениями.
- Если одновременно `synchronous_commit = off` и используется распределение Зипфа, то возникает проблема с high-contention write.
- Для high-contention write есть ряд патчей, улучшающих ситуацию.
- Можно успешно бороться с high-contention write на уровне приложения, хотя это и костыли.
- Мы всё понимаем, и будем и дальше трудиться над улучшением ситуации.



# СПАСИБО !





**Thanks !**