

Performance tuning

part II. instance tuning

Владимир Пудовченко



Что такое тюнинг?

"диагностика неисправностей и поиск и устранение неэффективностей"

Потоки получения информации:

1. логи субд
2. семплирование ожиданий
3. общая системная информация об утилизации основных ресурсов
4. представления с оперативными и накопительными статистиками `pg_stat_*` (внутренние системные статистики postgres)
5. Профили запросов (top-SQL)

USE – Utilisation Saturation Errors
RED – Rate Errors Duration
LETS – Latency Errors Traffic Saturation (Google 4 Golden Signals)
STELA – (LETS+A) + Availability

performance tuning

сопровождение

- ☐ настройка резервирования
- ☐ обновление
- ☐ мониторинг

поиск неисправностей

- ☐ субд почему-то упала
- ☐ упала и не включается
- ☐ регулярно падает

Подход к настройке

Задача тюнинга СУБД заключается в сокращении времени на обработку запроса сервером И оптимизации объявленных показателей (TPS, время отклика).

SQL tuning

Проблемное место известно. Долго загружается какая-то страница приложения или долго строится отчёт.

Instance tuning

- Широкий(количественный) взгляд на работу сервера (TPS , latency и тд.)
- Семплирование ожиданий
- Нахождение узкого места в работе сервера

Hardware\OS tuning

- Ulimits, ФС, SWAP.
- параметры ядра ОС sysctl
- BIOS/UEFI



Доклады Михаила Жилина:

- **Лучшие практики по OS + PostgreSQL**
- **Ускоряем железо и OS под PostgreSQL**
- **Разгоняем железо и операционную систему на максимальную производительность. Бенчмаркаем на PostgreSQL**

Простой вариант – tuned

- Позволяет системным администраторам легко и динамично менять настройки ядра
- Не нужно вносить изменения в /etc/sysctl

```
#sudo apt install tune
```

```
student:~$ tuned-adm list | grep postgresql
- postgresql - Optimize for PostgreSQL server
```

```
grep -v "^s*[\#;]\|^s*$" /usr/lib/tuned/postgresql/tuned.conf
[main]
summary=Optimize for PostgreSQL server
include=throughput-performance
[cpu]
force_latency=1
[vm]
transparent_hugepages=never
[sysctl]
vm.dirty_background_ratio = 0
vm.dirty_ratio = 0
vm.dirty_background_bytes = 67108864
vm.dirty_bytes = 536870912
vm.swappiness=3
kernel.sched_autogroup_enabled = 0
[scheduler]
sched_min_granularity_ns = 10000000
sched_migration_cost_ns = 50000000
```

Instance tuning

Instance tuning

цель: выжить из дорогостоящего железа максимум, обеспечить стабильную работу (без скачков).

1. первичная настройка экземпляра из коробки.
2. донастройка на основе мониторинга под специфическую нагрузку конкретного приложения.

В PG порядка 150+ GUC-параметров (postgresql.conf)

из которых 20-30 параметров для которых известны «правильные значения»

Параметры хранения (табличные параметры)

1. switchover узлов кластера под управлением кластерного ПО patroni из-за высокой загрузки процессора.
2. отставание синхронной реплики
3. вымывание кэша после ночных джобов
4. autovacuum упирается в потолок производительности -> разрастание таблиц и индексов

«медленно работает база»

ПОДХОДЫ

Повышение производительности – устранение проблем с производительностью

1. Для диагностики проблем с производительностью требуется знание внутреннего устройства СУБД:

- архитектура
- фоновые процессы
- области памяти и тд.

2. Книги рецептов не существует:

- каждый случай уникален
- 100% однозначных закономерностей не существует
- пока ещё нет автономных БД

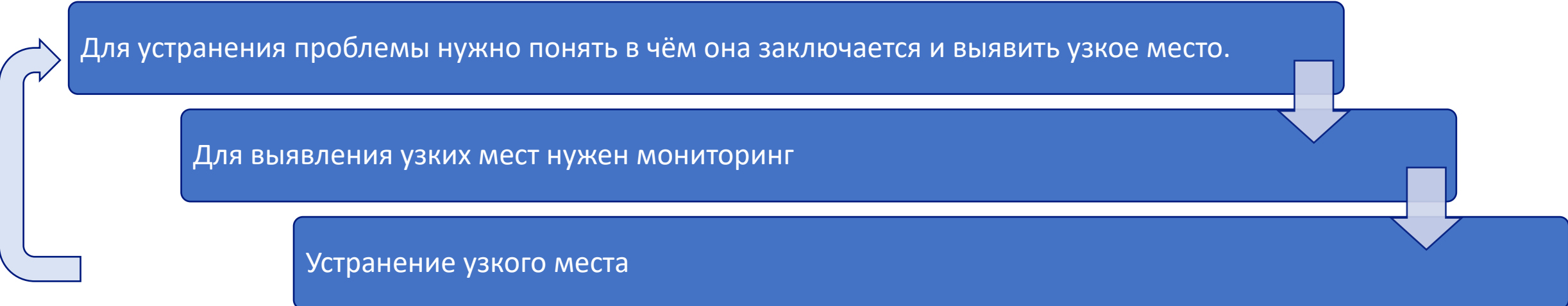


куда теряется производительность

Instance tuning

Количественный взгляд на работу СУБД (TPS, latency, etc.)

- растёт БД (как в ширину, так и в глубину)
появляется необходимость дополнительно что-то сохранять в БД
растут исторические данные
- увеличивается число пользователей
- появляется аналитика и тяжёлые OLAP-запросы



```

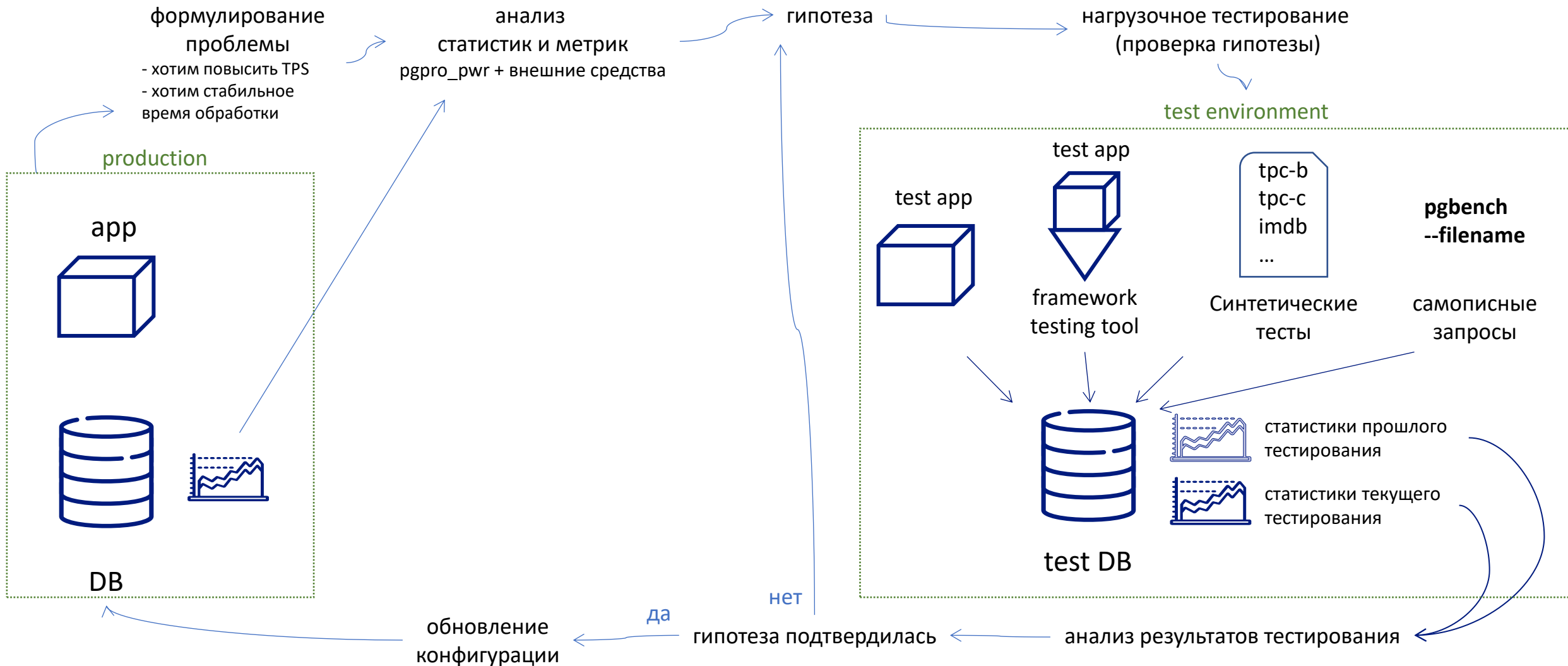
graph TD
    A[Для устранения проблемы нужно понять в чём она заключается и выявить узкое место.] --> B[Для выявления узких мест нужен мониторинг]
    B --> C[Устранение узкого места]
    C --> A
  
```

Для устранения проблемы нужно понять в чём она заключается и выявить узкое место.

Для выявления узких мест нужен мониторинг

Устранение узкого места

Тестирование (CI/CD pipeline)



Мониторинг внешними средствами

htop — информативная версия **top**.
atop — продвинутый монитор процессов
iotop — **мониторинг** ввода/вывода
И тд.

Если это выделенный под СУБД сервер, то и нагрузку на систему создаёт только СУБД.
Внешними средствами мы везде будем видеть только Postgres.

Внутренний мониторинг

`pg_stat_activity`. Мониторинг активностей и ожиданий

Когда процесс ожидает чего-либо, этот факт отражается в представлении **`pg_stat_activity`**

`wait_event_type` — тип ожидания

`wait_event` — имя конкретного ожидания

Информация только на текущий момент

единственный способ получить картину во времени — семплирование
достоверная картина только при большом числе измерений



Семплирование

pg_wait_sampling. Сбор статистик с некоторой периодичностью

pg_wait_sampling — это расширение Postgres Pro для периодического сбора статистики по событиям ожидания.

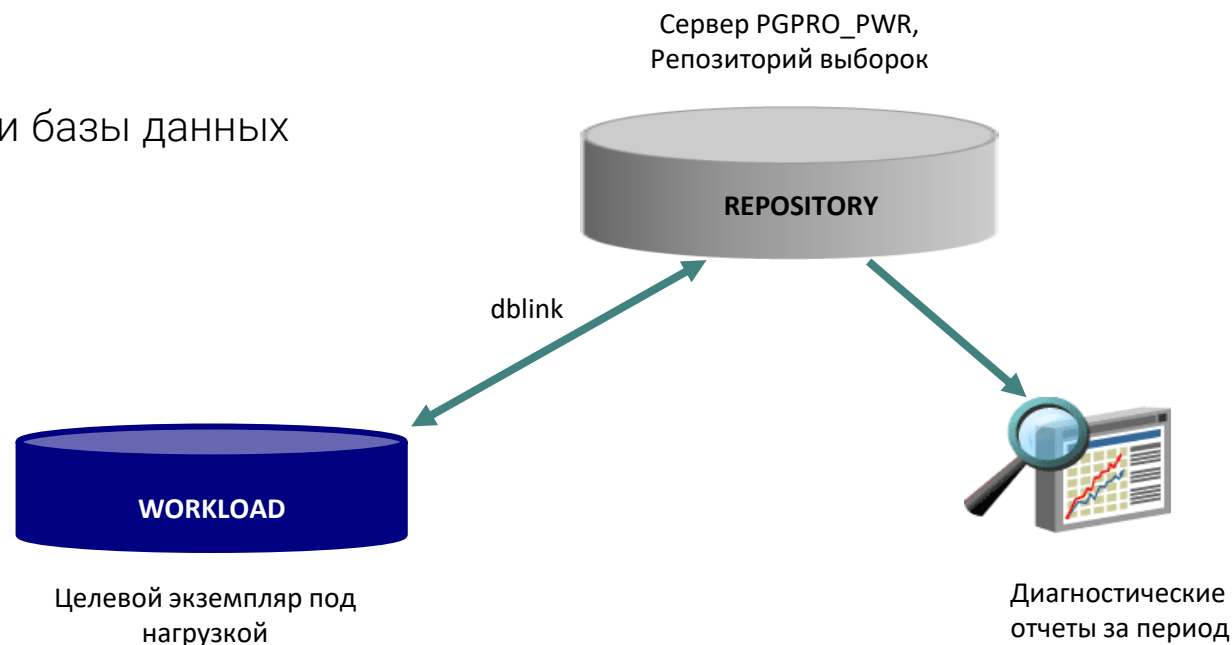
pid	backend_type	app	event_type	event	count
23559	checkpointer		Activity	CheckpointerMain	684
23559	checkpointer		IO	DataFileSync	5
23559	checkpointer		IO	SLRUFlushSync	4
23559	checkpointer		IO	ControlFileSyncUpdate	3
23559	checkpointer		IO	WALSync	2
23560	background writer		Activity	BgWriterMain	708
23561	walwriter		Activity	WalWriterMain	705
23561	walwriter		IO	WALSync	3
23562	autovacuum launcher		Activity	AutoVacuumMain	708
23565	logical replication launcher		Activity	LogicalLauncherMain	708
23600			IPC	CheckpointDone	24
23600			Client	ClientRead	10
23600			IO	CopyFileRead	1
23651	client backend	psql	Client	ClientRead	302
23651	client backend	psql	IO	WALSync	1
23621	client backend	psql	IO	WALSync	1
23621	client backend	psql	Activity	ClientRead	302

Отчёты pgpro_pwr

Абсолютные значения ожиданий сами по себе мало информативны. Нужно понимать за какой период они собирались. (семплирование для семплирования).

В каждый снимок pgpro_pwr попадает кумулятивная статистика ожиданий, после чего вызывается pg_wait_sampling_reset(). Построив отчёт по снимкам pgpro_pwr можно получить детальную картину о том какие ожидания преимущественно были между интервалов снимков.

- «Диагностическая карта» производительности базы данных
- Периодические снимки статистики - samples
- Репозиторий – хранилище истории снимков
- Детальный отчет за выбранный период



PGPRO_STATS – графический интерфейс PPEM



Экземпляры / PgPro Manager repository / SQL статистика

Статистика выражения

Query ID	Top level	База данных	Пользователь	Calls	Rows	Total execution time, mc	Total planning time, mc	Blocks time, mc
Plan ID						Max Min Mean Stdev	Max Min Mean Stdev	Read Write
2948794686727685322	-1764185593907006535	pgpro_manager_repo	pgpro_manager	12074	13229470	8568477.46	0.00	
						2059.56 27731 709.66 346.14	0.00 0.00 0.00 0.00	2712993.95 31215
-4940975716380762037	-7452334651223299015	pgpro_manager_repo	pgpro_manager	4512	4506	4658594.93	0.00	
						291213 1.46 1032.49 744.84	0.00 0.00 0.00 0.00	0.00 0.00
5977927359041152324	-1765265332860400087	pgpro_manager_repo	pgpro_manager	14022	14022	4411282.95	0.00	
						2681.74 110.75 314.60 141.39	0.00 0.00 0.00 0.00	100124.05 1162.03
-8345031607887110380	325442572100685824	pgpro_manager_repo	pgpro_manager	4512	4506	2319632.80	0.00	
						1512.71 1.01 514.10 371.05	0.00 0.00 0.00 0.00	0.00 0.00
-9009610577462668850	6701048979024298807	pgpro_manager_repo	pgpro_manager	4512	4512	2315056.72	0.00	
						1439.80 0.98 513.09 369.99	0.00 0.00 0.00 0.00	0.00 0.00
4049638393080033471	-2873485088126847906	pgpro_manager_repo	pgpro_manager	4512	4512	2283759.93	0.00	
						1424.92 1.00 506.15 365.18	0.00 0.00 0.00 0.00	0.00 0.00
-1532875410990978629	-8289805283452537256	pgpro_manager_repo	postgres	110465	3048834	1456303.19	0.00	
						59.68 0.07 1318 15.71	0.00 0.00 0.00 0.00	0.92 0.00
6962627314752790933	-74988376158504123	pgpro_manager_repo	pgpro_manager	69872847	69872847	1407035.74	0.00	
						4813 0.01 0.02 0.02	0.00 0.00 0.00 0.00	124.96 1.08

Экземпляры / PgPro Manager repository / stat_statements / Оператор

Query ID 2948794686727685322

Common	Execution time, mc	WAL	Shared blocks
Calls 12074	Total 8568477.46	Bytes, B 2174586882.00	Hit 1067149726
Rows 13229470	Max 2059.56	Records 13288957	Read 144314447
Пользователь pgpro_manager	Min 27731	Fpi 228513	Dirtied 237299
База данных pgpro_manager_repo	Mean 709.66		Written 4438
Query ID 2948794686727685322	Stdev 346.14	JIT	Local blocks
Plan ID -1764185593907006535	Rusage (63076304,21659648,8396.62913300003111,9.66044300000014,818831868.0,0.0,0.1832,572801)	Functions 0	Hit 0
Plans 0	Planning time, mc	Generation time, mc 0.00	Read 0
Top level true	Total 0.00	Optimization 0	Dirtied 0
	Max 0.00	Optimization time, mc 0.00	Written 0
	Min 0.00	Emission 0	Temp blocks
	Mean 0.00	Emission time, mc 0.00	Read count 0
	Stdev 0.00	Inlining 0	Read time, mc 0.00
	Rusage (0,0,0,0,0,0,0,0,0,0)	Inlining time, mc 0.00	Written count 0
	Blocks time, mc	Invalid msgs (0,0,0,0,0,0,0,0)	Written time, mc 0.00

Wait stats, mc

IO	
DataFileRead	10
DataFileWrite	10
LWLock	
BufferContent	30
Total	
IO	20
LWLock	30
Total	50

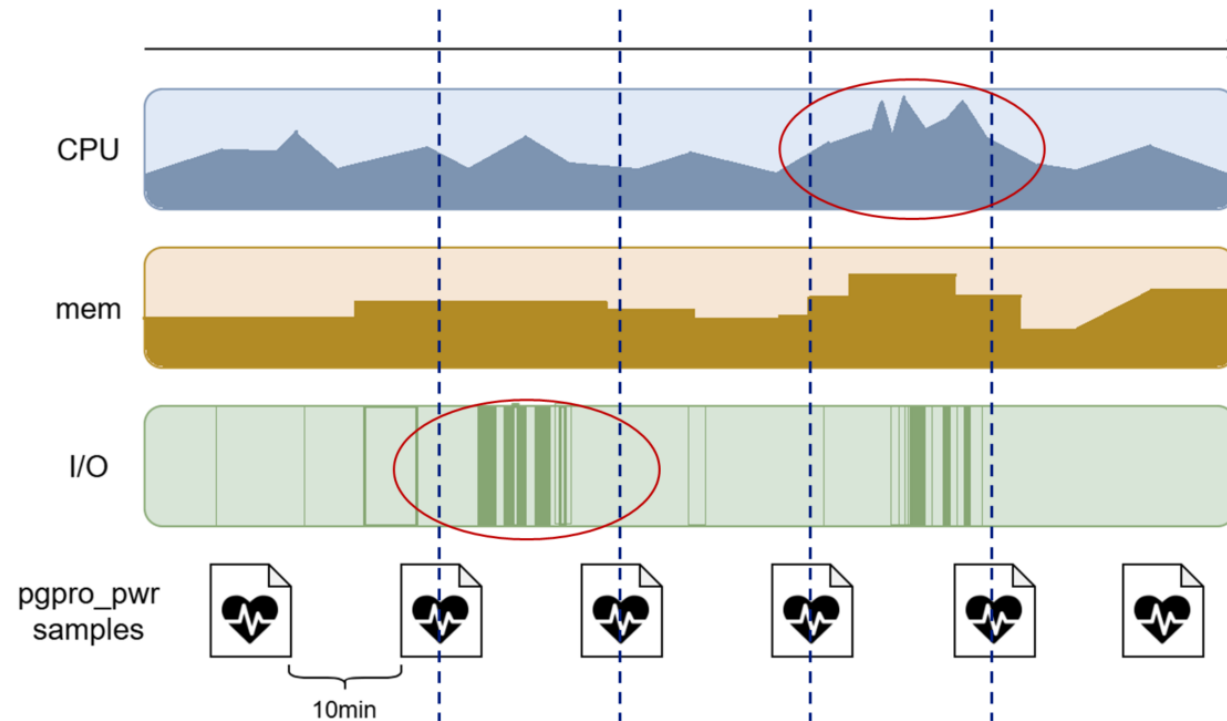
Query

```
DELETE FROM monitoring_history
WHERE
  value_set < (
    SELECT
      to_timestamp(
        extract(
          $1
        )
      FROM
        current_timestamp
      ) - $2::int
    )
  AND instance_id = $3
```

Plan

```
Delete on public.monitoring_history
InitPlan 1 (returns $0)
-> Result
  Output: to_timestamp(((EXTRACT($1 FROM CURRENT_TIMESTAMP) - $2))::double precision)
  -> Seq Scan on public.monitoring_history
    Output: monitoring_history.ctid
    Filter: ((monitoring_history.value_set < $0) AND (monitoring_history.instance_id = $3))
```


Дифференциальные отчёты pgpro_pwr



Отчёты pgpro_pwr в RPEM

← НАЗАД

ПРОФИЙЛЕР

Обзор

Построение отчетов

Отчеты

Снимки статистики

Расписание

Серверы

Экземпляры / PPro-ent-16 / Профайлер

Активность экземпляра

БД postgres

Сервер local

Период 05.11.2024, 00:00 - 06.11.2024, 00:00

ВЫБРАТЬ

PostgreSQL instance: tuples



PostgreSQL bgwriter buffers



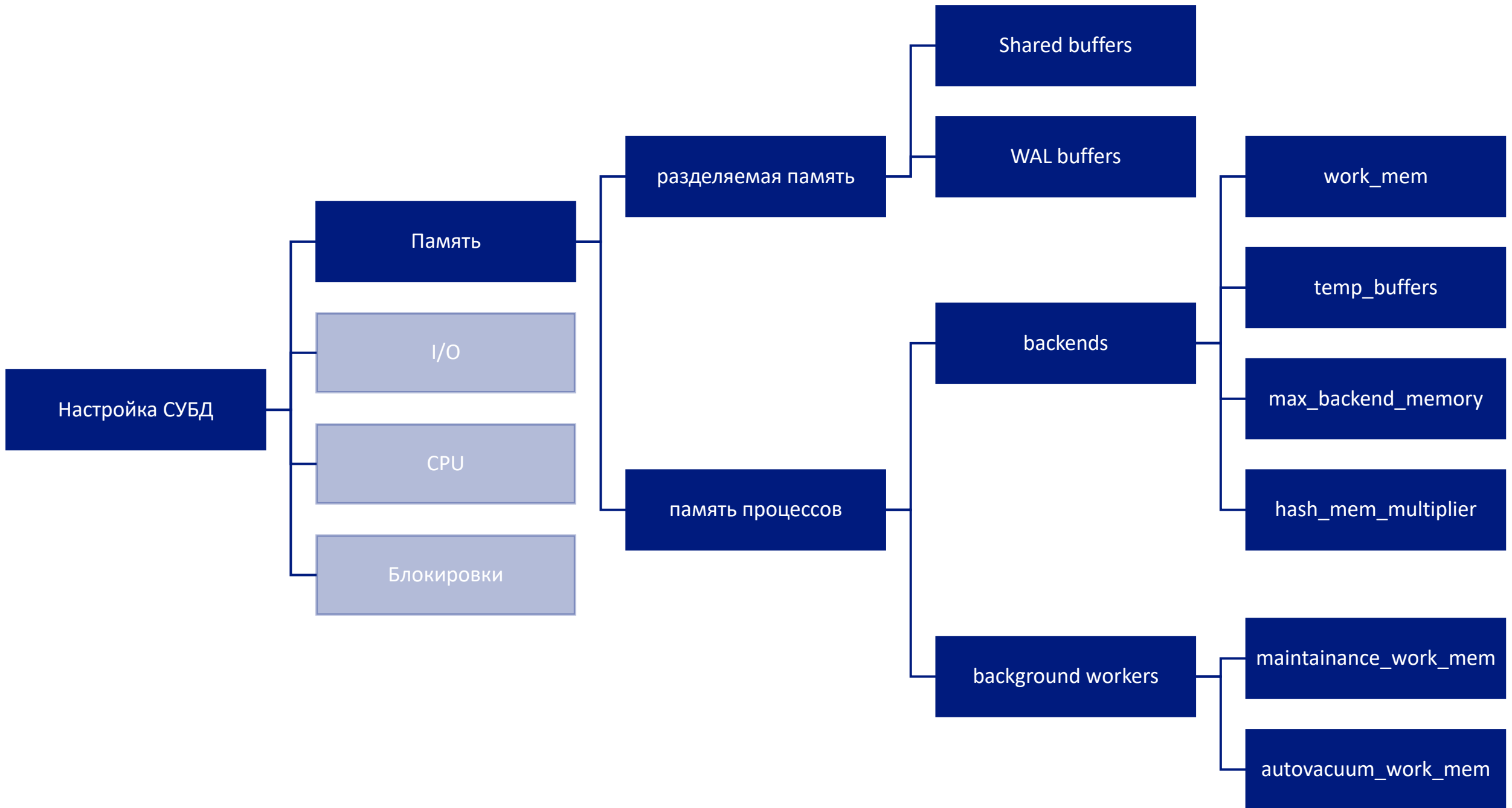
LOAD DISTRIBUTION AMONG HEAVILY LOADED USERS

Resource	Load distribution
Total time (sec.)	<div></div>
Executed count	<div></div>
I/O time (sec.)	<div></div>
Blocks fetched	<div></div>
Shared blocks read	<div></div>
Shared blocks dirtied	<div></div>
Shared blocks written	<div></div>
WAL generated	<div></div>
Temp and Local blocks written	<div></div>
Temp and Local blocks read	<div></div>
Invalidation messages sent	<div></div>
Cache resets	

SQL query statistics

TOP SQL BY EXECUTION TIME

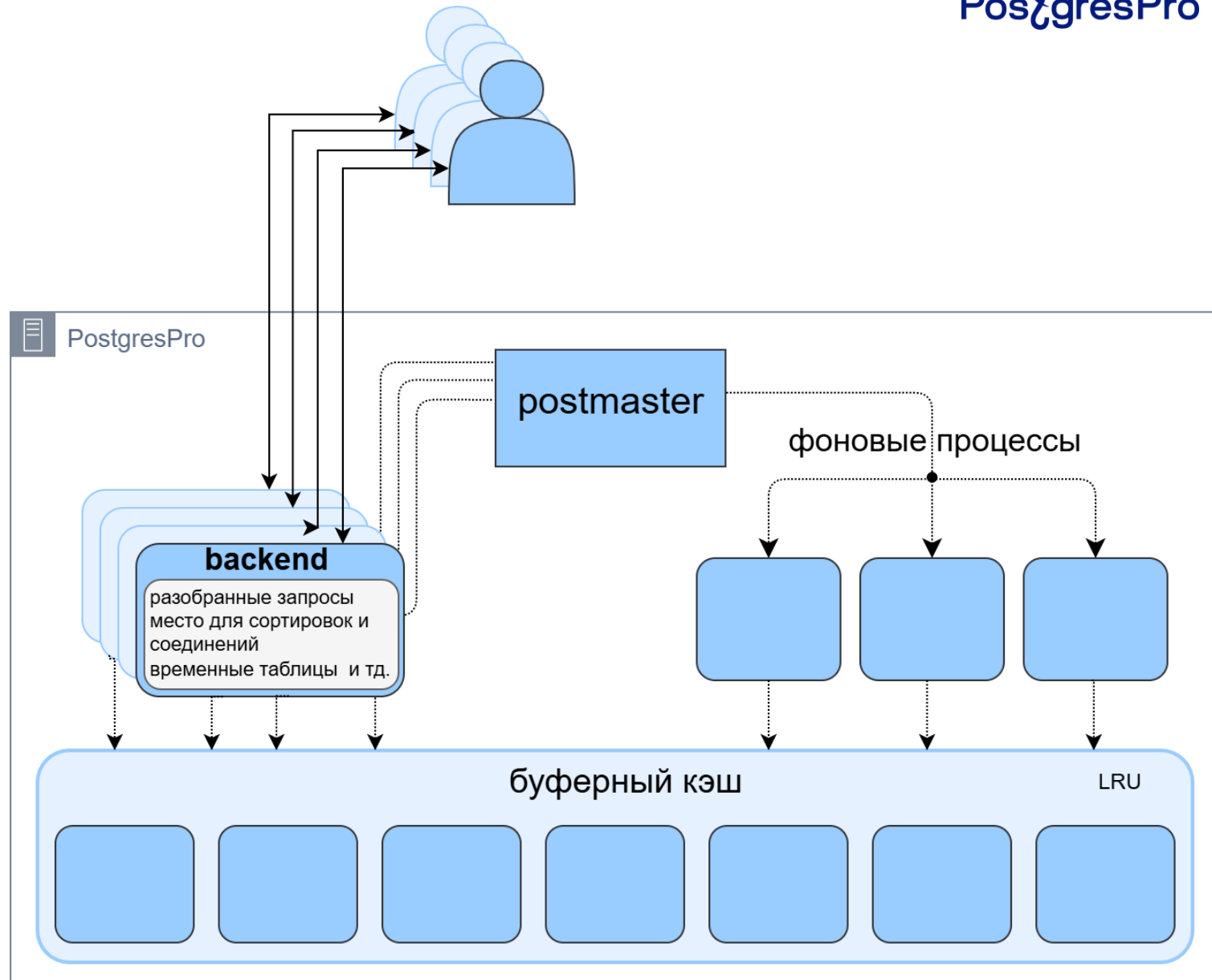
Query ID	Plan ID	Database	User	Exec (s)	%Total	I/O time (s)		CPU time (s)		Rows	Execution times (ms)				Executions
						Read	Write	Usr	Sys		Mean	Min	Max	StdErr	



Буферный кэш

В заголовке буфера:

- bufferid
- relfilenode
- reltablespace
- reldatabase
- relforknumber
- isdirty
- **usagecount**
- pinning_backends



Настройка shared_buffers

shared_buffers используется для определения того, сколько памяти будет выделено серверу для кэширования данных.

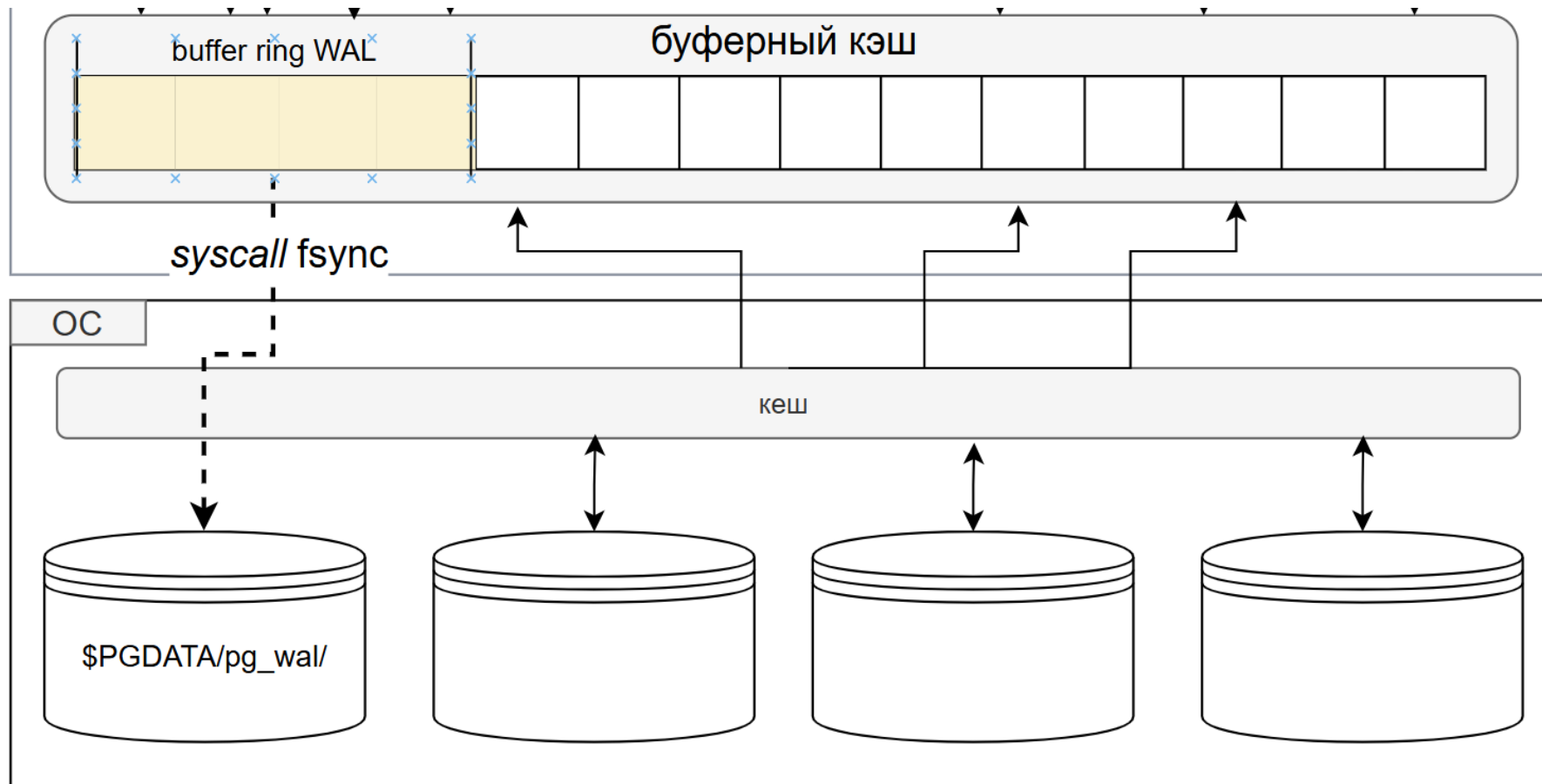
Разумное начальное значение - $\frac{1}{4}$ памяти на сервере. (99% случаев верно)

Для более прицельной настройки можно воспользоваться расширением **pg_buffercache** (в составе субд). Либо если серверные процессы часто ожидают блокировки «buffer pin lock»

Инспекция кэша СУБД:

```
postgres=# select usagecount, count(*) from pg_buffercache group by usagecount order by usagecount ;
 usagecount | count
-----+-----
          1 |    1029
          2 |     517
          3 |      33
          4 |     673
          5 |    6544
    <NULL> | 245156
(6 строк)
```

wal_buffers - объём разделяемой памяти, который будет использоваться для буферизации данных WAL, ещё не записанных на диск. Значение по умолчанию, равное -1, задаёт размер, равный 1/32 (около 3%) от shared_buffers



Мониторинг wal_buffers

Стоит увеличивать если:

1. Серверные процессы часто ожидают получения блокировок «WALinsertLock» либо «WALBufMappingLock».
2. Реплика периодически отстаёт, тогда процессу **wal sender** приходится читать недостающие записи с диска.

Это можно определить:

- Процесс wal sender имеет не нулевой i/o на мастере. Большая разница sent_lsn и write_lsn

Использование временных таблиц

temp_buffers – ограничивает максимальный объем памяти для хранения данных временных таблиц (в рамках каждого отдельного сеанса).

CREATE { TEMPORARY | TEMP } TABLE

Как определить, что используются временные таблицы?

Отчет pg_profile или pgpro_pwr!

Statement statistics by database

Database	Calls	Time (s)				Fetched (blk)		Dirtied (blk)		Temp (blk)		Local (blk)	
		Exec	Read	Write	Trg	Shared	Local	Shared	Local	Read	Write	Read	Write
	21142063357	5434617.12	13951.10	55.84		3675989789344	828719501	1039147932	149806370	294562125	305989125	710679390	326520283
pg_profile	24942263	7769.07	0.74			945853344	94006	4794711	6458	5329774	7327746	7074	9846
postgres	1983052	3119.95	0.01			60947665		3					
template1	24	0.00				67							
Total	21168988696	5445506.14	13951.85	55.84		3676996590420	828813507	1043942646	149812828	299891899	313316871	710686464	326530129

maintenance_work_mem задаёт максимальный объём памяти для операций обслуживания БД, в частности VACUUM, CREATE INDEX и ALTER TABLE ADD FOREIGN KEY. Увеличение этого значения может привести к ускорению операций очистки и восстановления БД из копии.

autovacuum_work_mem задаёт максимальный объём памяти, который будет использовать каждый рабочий процесс автоочистки. (при -1 используется maintenance_work_mem).

- ! Выделяется в полном объёме и может выделяться autovacuum_max_workers раз, поэтому не стоит устанавливать значение по умолчанию слишком большим.
- ! Недостаточный размер приводит к многократному сканированию индексов таблицы при очистке.

Мониторинг autovacuum_work_mem

Число индексных сканирований должно равняться числу индексов по таблице.

- `pg_stat_progress_vacuum.index_vacuum_count`
- `VACUUM (verbose) <tablename>`

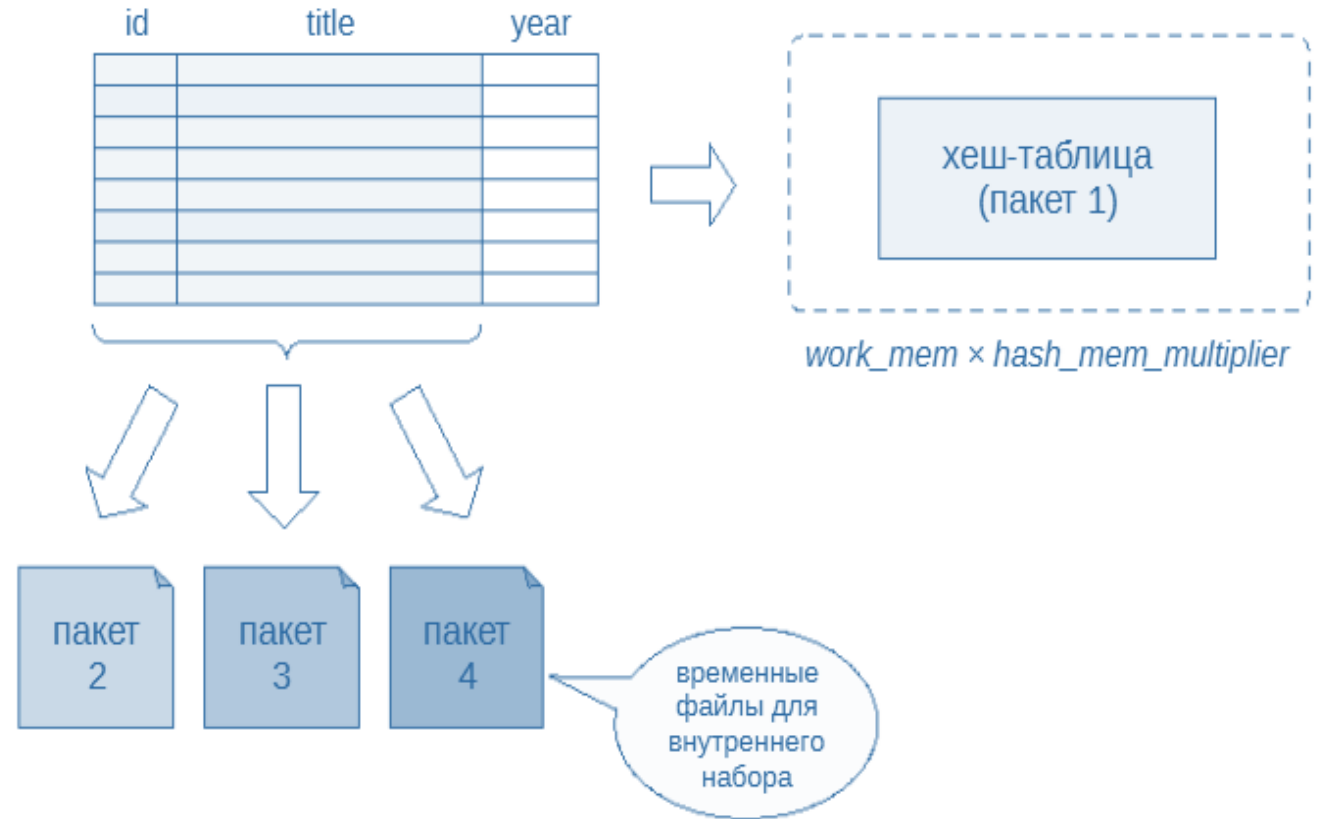
```
postgres=# vacuum (verbose) t ;  
INFO:  vacuuming "postgres.public.t"  
INFO:  finished vacuuming "postgres.public.t": index scans: 1  
pages: 0 removed, 6431 remain, 2054 scanned (31.94% of total)  
tuples: 189999 removed, 870865 remain, 0 are dead but not yet r
```

- `log_autovacuum_min_duration TO '1ms'`

>>	Дата, время	Экземпляр	Пользователь	База данных ▼	Тип ▼
^	5.11.2024, 21:38:03.63	PPro-ent-16			LOG
automatic vacuum of table "postgres.profile.last_stat_kcache_srv1": index scans: 1pages: 3 removed, 2 re cutoff: 7150, which was 1 XIDs old when operation endednew relfrozenxid: 7158, which is 16 XIDs ahead of of total) had 0 tuples frozenindex scan needed: 4 pages from table (80.00% of total) had 59 dead item id					

work_mem - Задаёт базовый максимальный объём памяти, который будет использоваться во внутренних операциях при обработке запросов (например, для сортировки или хеш-таблиц), прежде чем будут задействованы временные файлы на диске.

hash_mem_multiplier - Используется для определения максимального объёма памяти, который может выделяться для операций с хешированием. Итоговый объём определяется произведением `work_mem` и `hash_mem_multiplier`.



Мониторинг work_mem

Как определить, что объем work_mem достаточен?

1. Отчет pg_profile или pgpro_pwr!

Statement statistics by database

Database	Calls	Time (s)				Fetched (blk)		Dirtied (blk)		Temp (blk)		Local (blk)	
		Exec	Read	Write	Trg	Shared	Local	Shared	Local	Read	Write	Read	Write
	21142063357	5434617.12	13951.10	55.84		3675989789344	828719501	1039147932	149806370	294562125	305989125	710679390	326520283
pg_profile	24942263	7769.07	0.74			945853344	94006	4794711	6458	5329774	7327746	7074	9846
postgres	1983052	3119.95	0.01			60947665		3					
template1	24	0.00				67							
Total	21168988696	5445506.14	13951.85	55.84		3676996590420	828813507	1043942646	149812828	299891899	313316871	710686464	326530129

2. Параметр log_temp_files = '0'; + log_statements = 'ALL';

Позволяет выявить top-sql запросы, которые не поместились в work_mem, и для выполнения которых потребовалось создание временных файлов на диске.

3. PPEM метрика «Temp files created»

Сколько памяти нужно выделить?

shared_buffers +

(work_mem * hash_mem_multiplier + temp_buffers) * max_connections +

work_mem * max_worker_processes +

autovacuum_max_workers * autovacuum_work_mem

можно ограничить
максимальное использование
памяти **max_backend_memory**

- не должен превышать максимальный объем памяти на сервере, чтобы избежать принудительного завершения основного процесса PostgreSQL OOM killer. Можно защититься отдав пространство под SWAP.

max_connections – не нужно выставлять огромным, т.к. это резервирует память.

Прогрев кэша

Проблема: после перезагрузки СУБД кэш пустой.

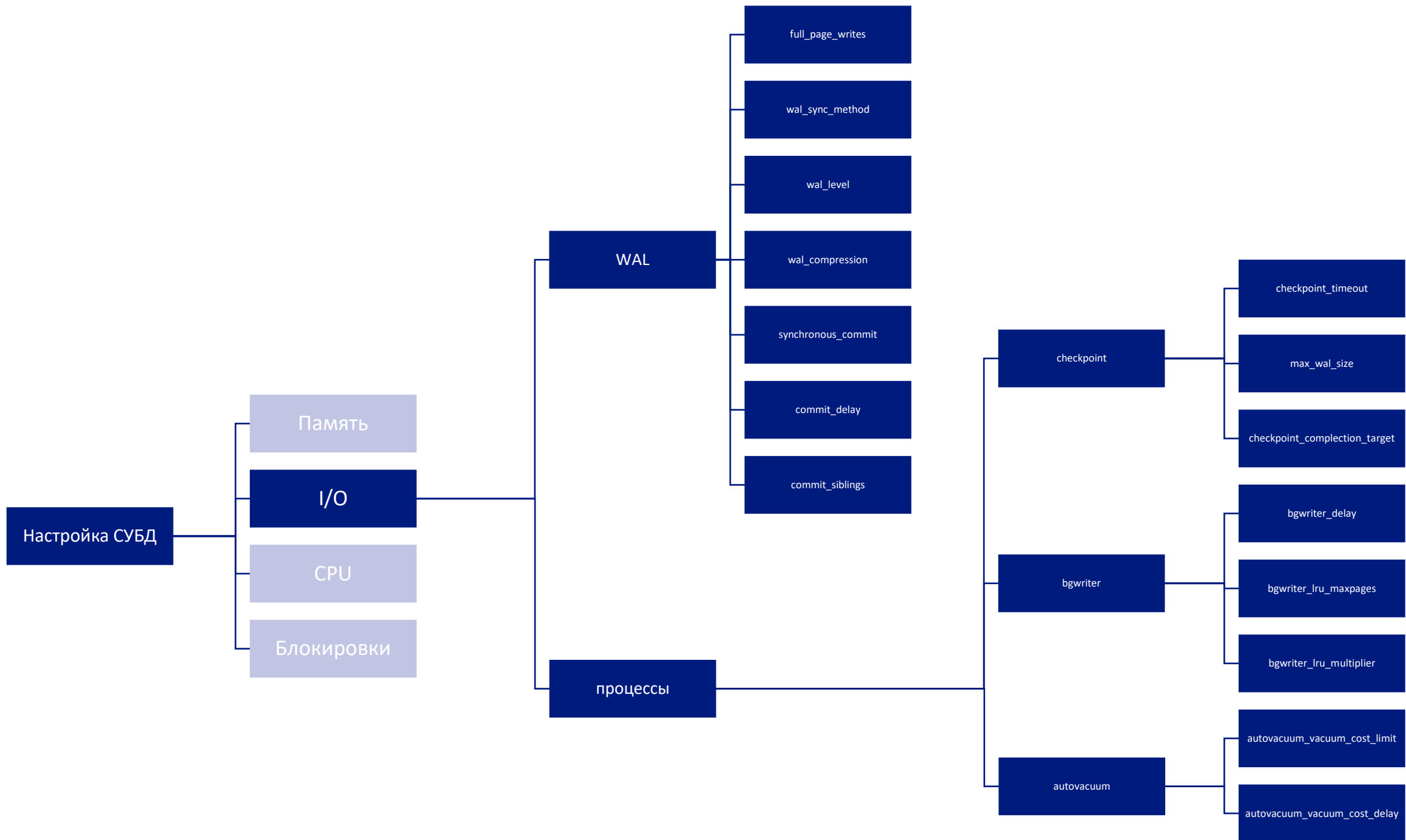
Решение: `pg_prewarm`

```
# postgresql.conf
```

```
shared_preload_libraries = 'pg_prewarm'
```

```
pg_prewarm.autoprewarm = true
```

```
pg_prewarm.autoprewarm_interval = 300s --Задаёт интервал между обновлениями файла autoprewarm.blocks.  
(mapping shared cache)
```



производительность

надёжность

✱ **~~fsync~~ = {on, off}**

wal_sync_method {

- open_datasync
- fdasync
- fsync
- open_sync
- ...

full_page_writes = {on, off}

Утилита **pg_test_fsync** позволяет выбрать способ, наиболее подходящий для конкретной ОС и конкретной файловой системы.

```
root@sa:~# /opt/pgpro/ent-17/bin/pg_test_fsync
на тест отводится 5 сек.
O_DIRECT на этой платформе не поддерживается для open_datasync и open_sync.
```

Сравнение методов синхронизации файлов при однократной записи 8 КБ:
(в порядке предпочтения для "wal_sync_method", за исключением того, что в Linux предпочитается fdasync)

open_datasync	106,654 оп/с	9376 мкс/оп
fdasync	103,019 оп/с	9707 мкс/оп
fsync	52,562 оп/с	19025 мкс/оп
fsync_writethrough	н/д	
open_sync	50,871 оп/с	19658 мкс/оп

Сравнение методов синхронизации файлов при двухкратной записи 8 КБ:
(в порядке предпочтения для "wal_sync_method", за исключением того, что в Linux предпочитается fdasync)

open_datasync	54,458 оп/с	18363 мкс/оп
fdasync	104,376 оп/с	9581 мкс/оп
fsync	52,102 оп/с	19193 мкс/оп
fsync_writethrough	н/д	
open_sync	26,865 оп/с	37223 мкс/оп

WAL



wal_compression = pglz,lz4,zstd,on,off - только для full_page_writes

synchronous_commit =

(можно менять на уровне транзакции)


remote_write
remote_apply
on
local
off

- Периоды синхронизации каждые **wal_writer_delay** (=200ms по умолчанию).
- Гарантируется согласованность, но не долговечность
- Зафиксированные транзакции могут пропасть за $(3 \times \text{wal_writer_delay})$

Где можно выиграть без потери надёжности:

commit_delay (ms)

commit_siblings (int)

Алгоритм: При числе одновременно работающих в системе транзакций более чем **commit_siblings**, процесс делает паузу, определяемую параметром **commit_delay**, в расчёте на то, что за время ожидания часть транзакций успеют завершиться и можно будет синхронизировать из записи за одну операцию.

Checkpoint

Периодический сброс всех грязных буферов на диск

- гарантирует попадание на диск всех изменений до контрольной точки
- ограничивает размер журнала, необходимого для восстановления

Частота срабатывания:

`checkpoint_timeout`

`max_wal_size`

Интенсивность работы:

`checkpoint_completion_target` = 0.9

14

Подход к настройке

Вытесняет на диск грязные страницы из буферного кэша в фоне, позволяя пользовательским процессам не тратить на это время.

<code>bgwriter_delay</code>	<code>= 200ms</code>
<code>bgwriter_lru_maxpages</code>	<code>= 100</code>
<code>bgwriter_lru_multiplier</code>	<code>= 2</code>

С маленькими значениями `bgwriter_lru_maxpages` и `bgwriter_lru_multiplier` уменьшается активность ввода/вывода со стороны процесса фоновой записи, но увеличивается вероятность того, что запись придётся производить непосредственно серверным процессам, что замедлит выполнение запросов.

Instance tuning

Сначала настраивается контрольная точка из соображений допустимого времени на восстановление, затем донастраивается bgWriter, таким образом, чтобы пользовательским процессам не приходилось вытеснять грязные буферы из кэша.

pg_stat_bgwriter. buffers_checkpoint

Количество буферов, записанных при выполнении контрольных точек

pg_stat_bgwriter. buffers_clean

Количество буферов, записанных фоновым процессом записи

pg_stat_bgwriter. buffers_backend

Количество буферов, записанных самим серверным процессом

15

pg_stat_io

`select writes,write_time from pg_stat_io where backend_type = 'client backend' and object = 'relation' ;`

VACUUM высвобождает пространство, занимаемое «мёртвыми» кортежами. При обычных операциях PostgreSQL кортежи, удалённые или устаревшие в результате обновления, физически не удаляются из таблицы; они сохраняются в ней, пока не будет выполнена команда VACUUM.

Регулирование агрессивности(производительности) автоматической очистки:

Процесс чередует работу и ожидание:

примерно vacuum_cost_limit условных единиц работы, затем засыпает на vacuum_cost_delay мс

```
vacuum_cost_limit = 200  
autovacuum_vacuum_cost_delay = 2ms
```

Определяет общий пул работы для всех
autovacuum_max_workers

bloat таблиц и индексов

Механизм MVCC предполагает появление мёртвых строк.

Негативное влияние:

- перерасход места на диске
- замедление последовательного просмотра таблиц
- уменьшение эффективности индексного доступа (лишние переходы из индекса на мёртвые табличные строки, т.к. в индексной записи нет информации о статусе табличной строки)

The diagram illustrates the lifecycle of data versions over time, represented by the x-axis (xid). The y-axis represents different versions of the data (ver_1, ver_2, ver_3, ver_4). Operations are marked as 'insert', 'update', or 'delete' at specific xid values.

- ver_1:** Starts with an 'insert' at xid 10. It continues as a solid red line until xid 110, where it transitions to a solid yellow line.
- ver_2:** Starts at xid 110 with an 'update' (yellow dot). It continues as a solid yellow line until xid 130, where it transitions to a solid dark blue line.
- ver_3:** Starts at xid 130 with an 'update' (dark blue dot). It continues as a solid dark blue line until xid 140, where it transitions to a solid yellow line.
- ver_4:** Starts at xid 140 with an 'update' (yellow dot). It continues as a solid yellow line.

Operations and transitions are marked with dots and labels:

- ver_1:** 'insert' at xid 10, 'update' at xid 110.
- ver_2:** 'update' at xid 110, 'update' at xid 130.
- ver_3:** 'update' at xid 130, 'update' at xid 140.
- ver_4:** 'update' at xid 140.

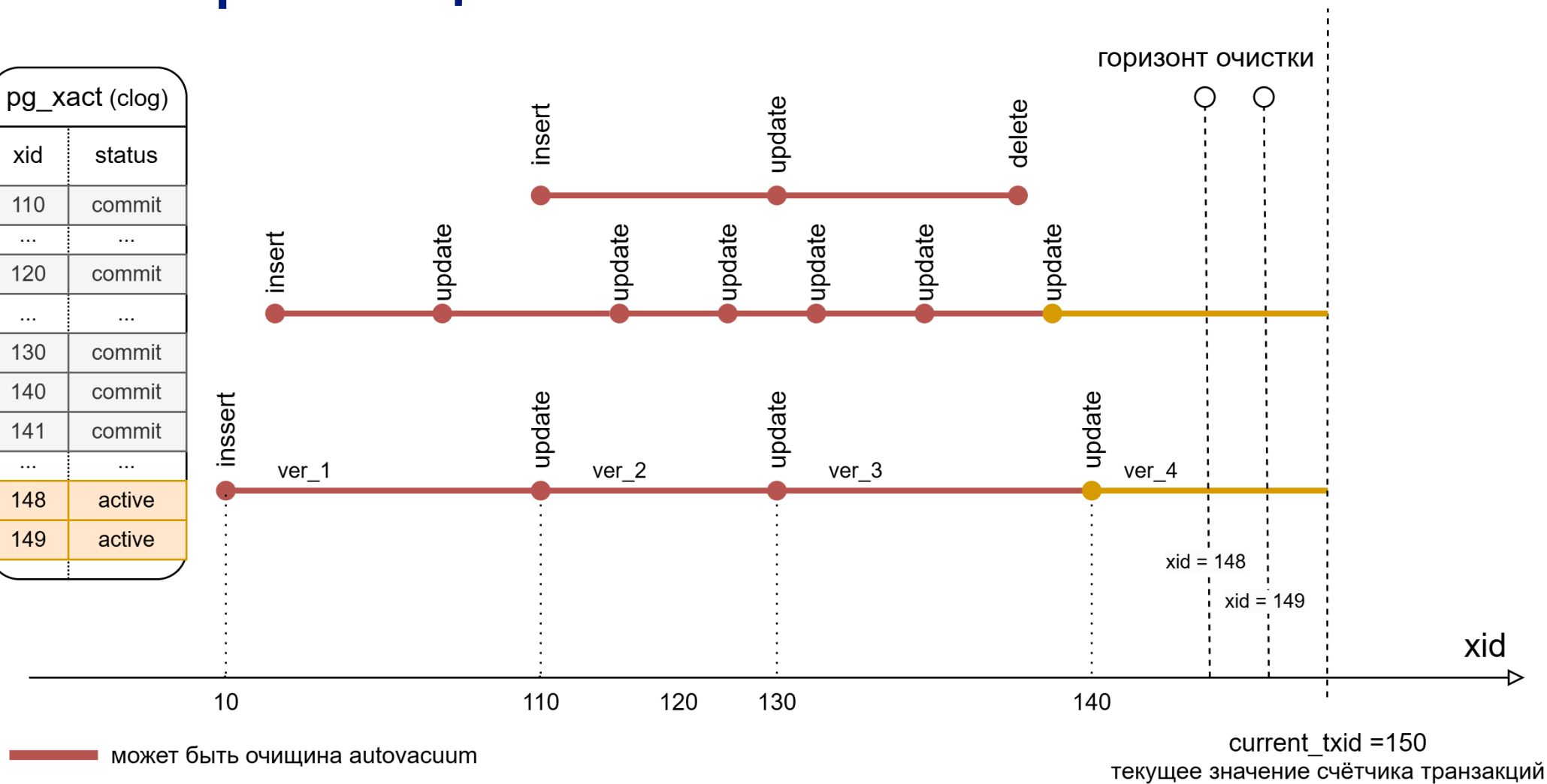
A vertical dashed line labeled 'горизонт очистки' (cleanup horizon) is positioned at xid 120. Other vertical dashed lines are at xid 148 and xid 149, marked with open circles. The x-axis is labeled 'xid' and has major ticks at 10, 110, 120, 130, and 140.

■ неактуальная версия, не видна ни в какой транзакции, но не может быть очищена

43

Долгие транзакции

pg_xact (clog)	
xid	status
110	commit
...	...
120	commit
...	...
130	commit
140	commit
141	commit
...	...
148	active
149	active



— может быть очищена autovacuum

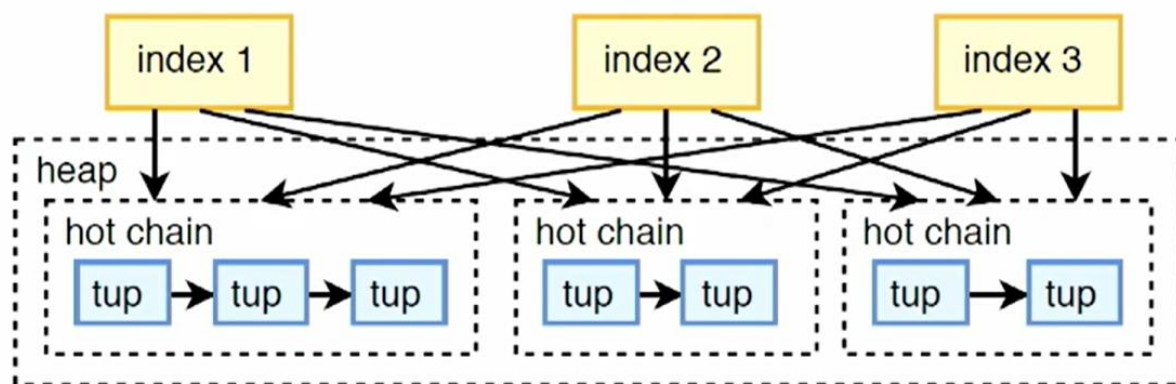
— попадает в снимок активной транзакции, не может быть очищена

— неактуальная версия, не видна ни в какой транзакции, но не может быть очищена

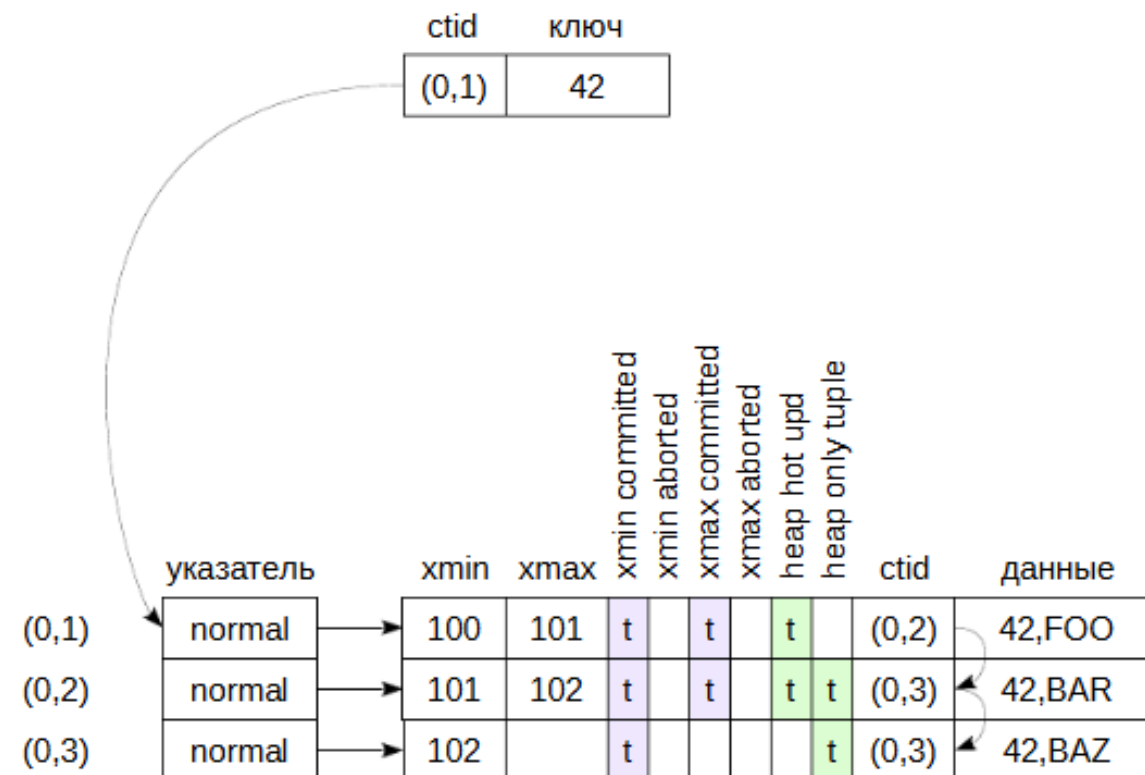
HOT-update

heap only tuple. Цепочка обновлений — только в пределах одной страницы

Для высоконагруженных UPDATE-ами таблиц.



- если в табличной странице не хватает места для новой версии, цепочка обрывается (как если бы оптимизация не работала)
- место в странице можно зарезервировать, уменьшив параметр хранения таблицы **fillfactor** (100 % → 10 %)

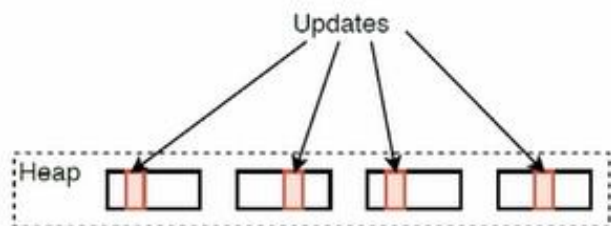


fillfactor

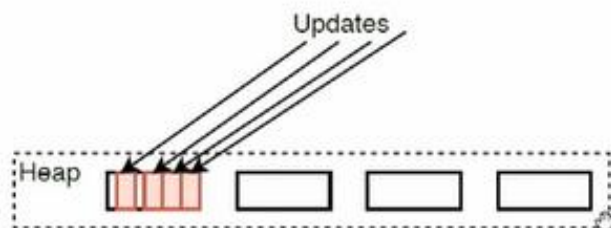
Плотность заполнения блоков

ALTER TABLE ... SET (fillfactor = 70);

Хорошо для HOT

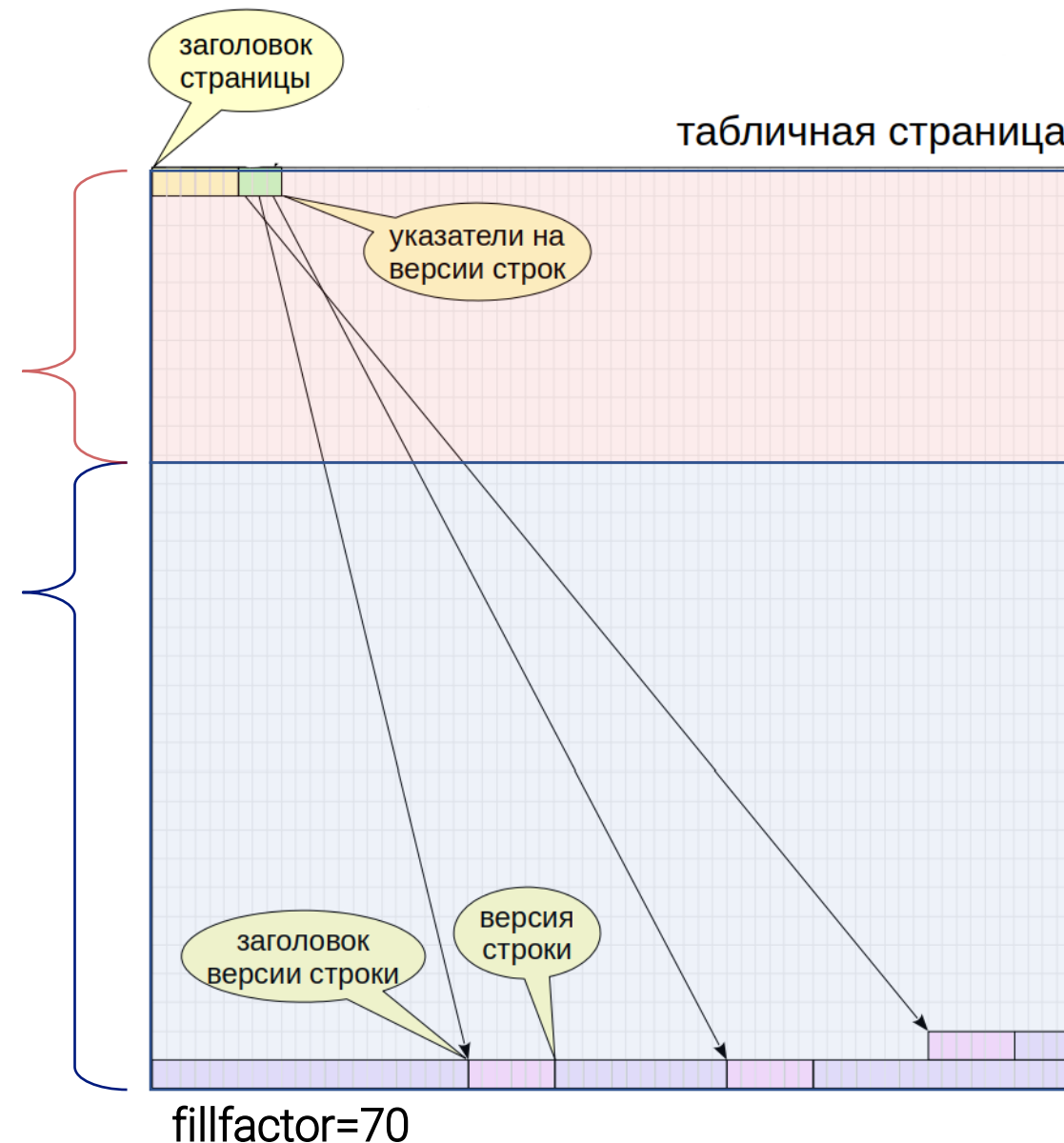


Плохо для HOT



Резерв под UPDATE

Свободное место



Настройка fillfactor

Большая разница между:

pg_stat_all_tables.n_tup_upd

pg_stat_all_tables.n_tup_hot_upd

=> Уменьшать fillfactor

```
postgres=# select relname, n_tup_upd, n_tup_hot_upd from pg_stat_all_tables where relname ~ 'pgbench_accounts' \gx
-[ RECORD 1 ]-+-----
relname      | pgbench_accounts
n_tup_upd    | 1245448
n_tup_hot_upd | 226944
```

Стандартные рекомендации по сокращению I/O

- **Выделенный СХД под WAL**

Заменить `$PGDATA/pg_wal` симлинком, указывающим на другое расположение

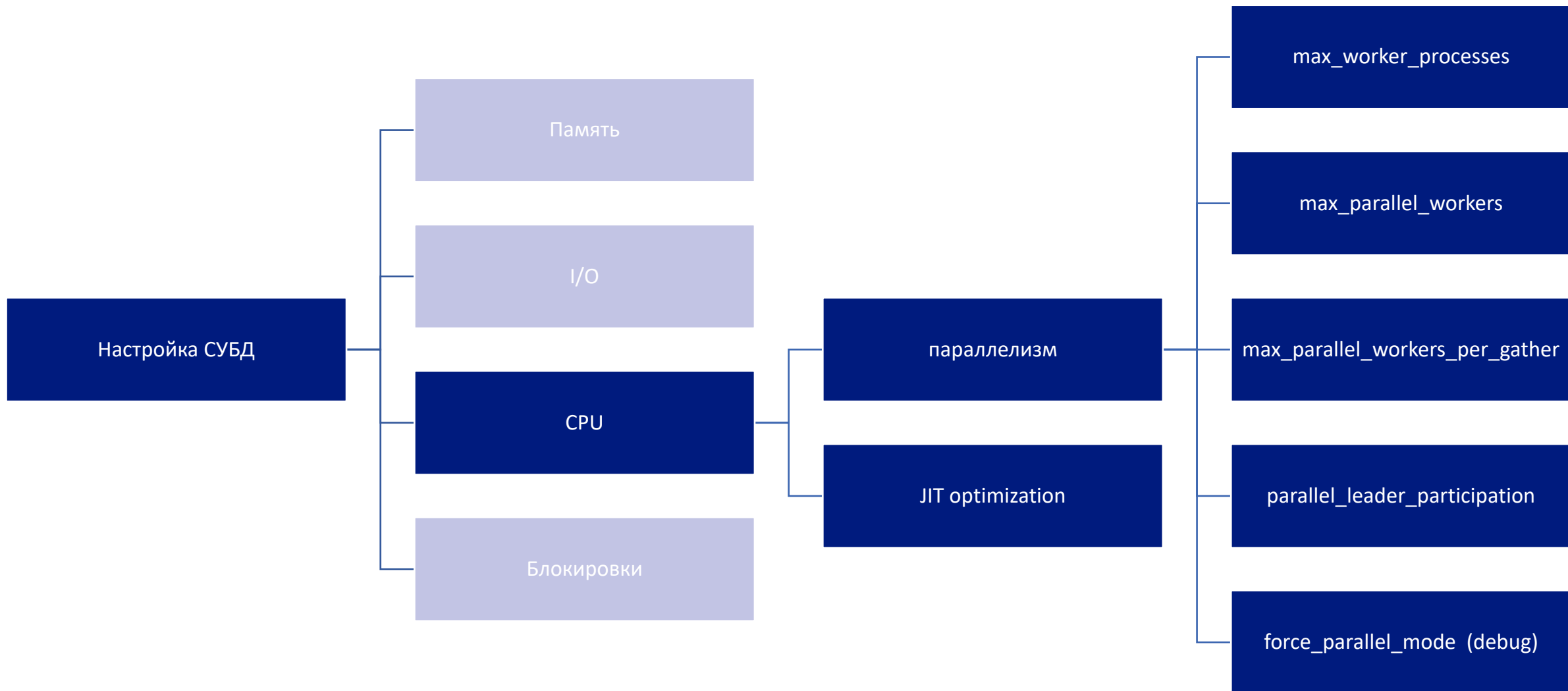
- **`temp_tablespace` (string)**

Эта переменная задаёт табличные пространства, в которых будут создаваться временные объекты (временные таблицы и индексы временных таблиц), которые не поместились в **`temp_buffers`**. В этих табличных пространствах также создаются временные файлы для внутреннего использования, если не поместились в **`work_mem`**.

- **Размещение `temp_tablespace` в tmpfs**

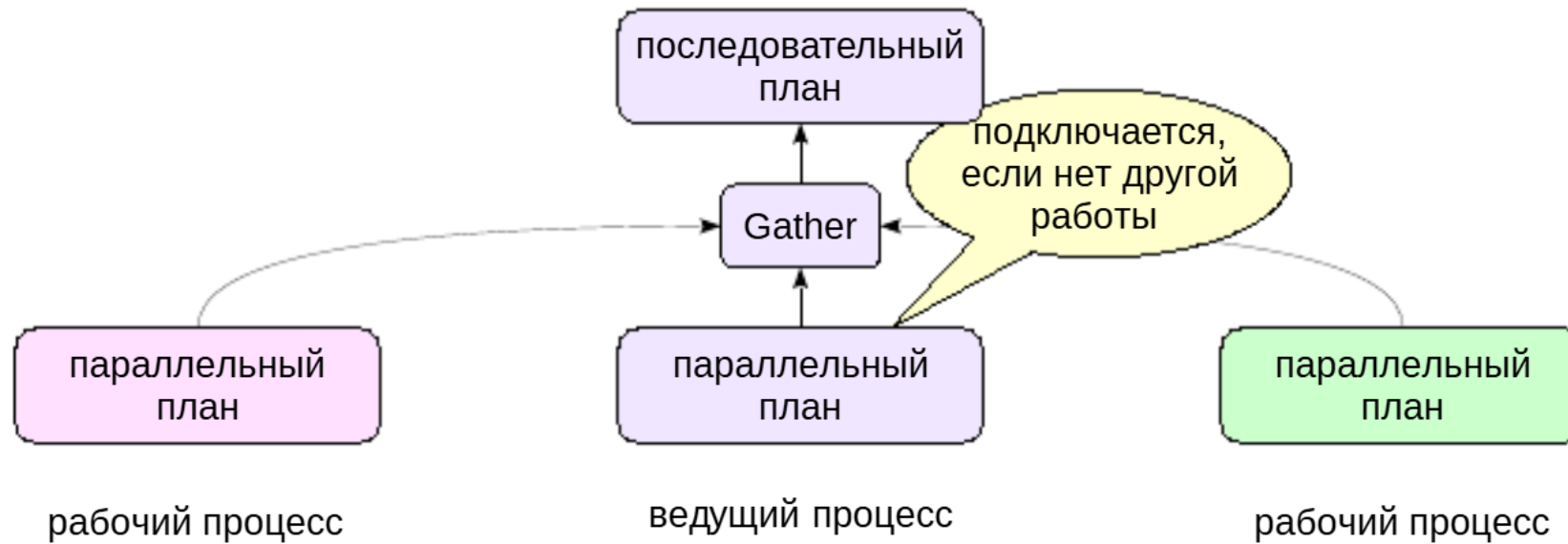
- **Секционирование нагруженных таблиц и размещение секций в разных пространствах (ILM).**

- **`CFS` (сжатый кэш на уровне ОС) + `wal_compression`**



Параллелизм

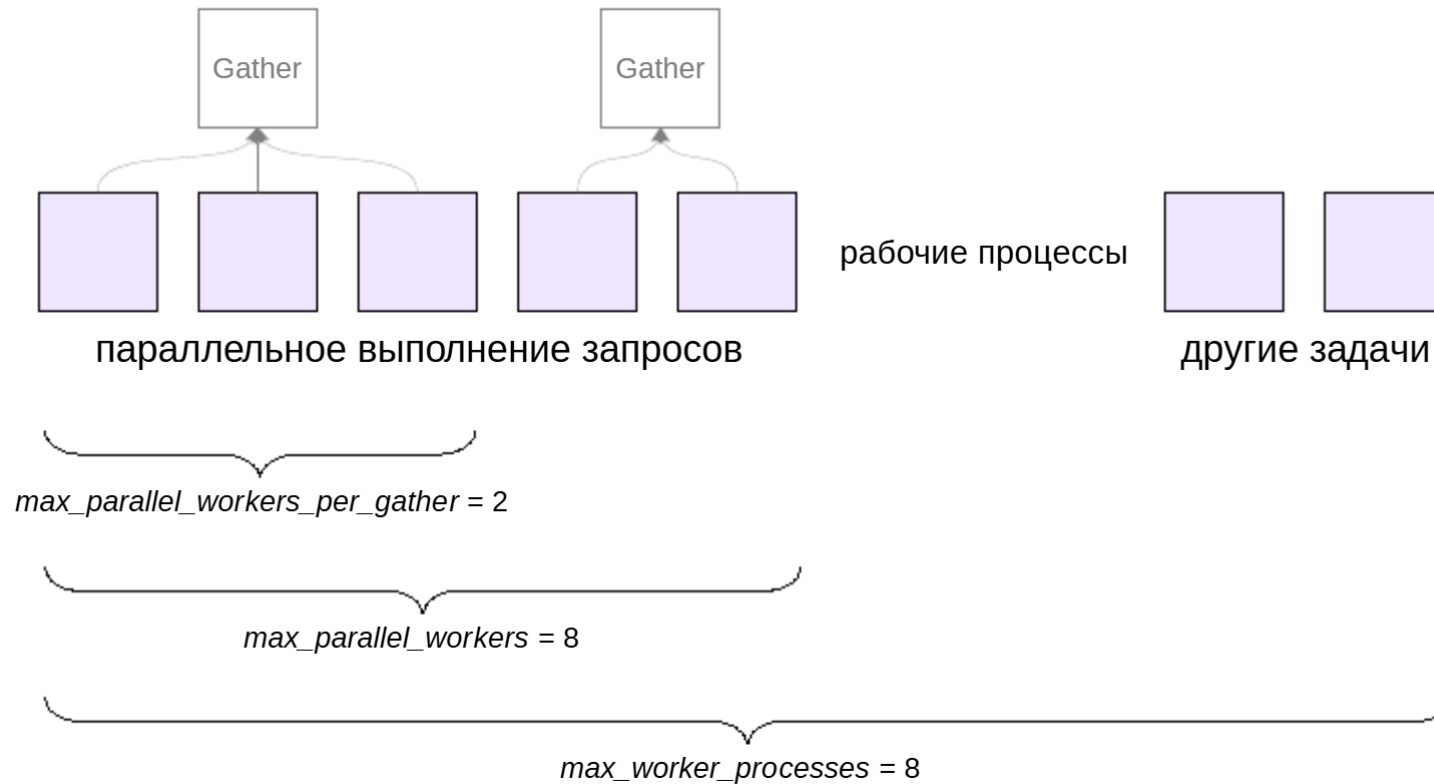
сжечь больше ресурсов, чтобы быстрее выполнить задачу



Параллелизм

Параметр конфигурации	Описание параметра	Значение по умолчанию
<code>max_worker_processes</code>	Задаёт максимальное число фоновых процессов, которое можно запустить в текущей системе.	8
<code>max_parallel_workers</code>	<p>Задаёт максимальное число рабочих процессов, которое система сможет поддерживать для параллельных операций.</p> <p>Заметьте, что значение данного параметра, превышающее <code>max_worker_processes</code>, не будет действовать, так как параллельные рабочие процессы берутся из пула рабочих процессов, ограничиваемого этим параметром</p>	8
<code>max_parallel_workers_per_gather</code>	Задаёт максимальное число рабочих процессов, которые могут запускаться одним узлом Gather	2

Параллелизм



Возникающие ожидания в параллельных планах

Событие ожидания **IPC**

Описание

MessageQueueInternal

Ожидание подключения другого процесса к общей очереди сообщений.

MessageQueuePutMessage

Ожидание записи сообщения протокола в общую очередь сообщений.

MessageQueueReceive

Ожидание получения байтов из общей очереди сообщений

MessageQueueSend

Ожидание передачи байтов в общую очередь сообщений

Накладные расходы

pid	wait_event	wait_event_type	duration	state	query
5930	[null]	[null]	00:00:03.452916	active	select subject.subject_id, addr.address_id.....
8665	MessageQueueSend	IPC	00:00:03.452310	active	select subject.subject_id, addr.address_id...
11330	MessageQueueSend	IPC	00:00:03.322182	active	select subject.subject_id, addr.address_id.....
13185	[null]	[null]	00:00:02.452919	active	select subject.subject_id, addr.address_id.....
17007	MessageQueueSend	IPC	00:00:02.442918	active	select subject.subject_id, addr.address_id.....
31220	MessageQueueSend	IPC	00:00:02.442917	active	select subject.subject_id, addr.address_id.....
47051	[null]	[null]	00:00:02.432916	active	select subject.subject_id, addr.address_id.....
55207	MessageQueueSend	IPC	00:00:02.432915	active	select subject.subject_id, addr.address_id.....

JIT

Just-In-Time compilation, компиляция «во время»

- JIT компиляция выполняется непосредственно во время выполнения запроса.
- JIT компиляция преобразует интерпретируемый вариант исполнения программы в машинный код.
- Применяется для ускорения вычисления сложных выражений в предикатах запросов.
- Если сложную функцию, вычисляемую в предикате, скомпилировать в машинный код, который выполняет процессор, то она, возможно, будет выполняться быстрее.

Оптимизации:

- Преобразование кортежей (Tuple deforming) — перевод кортежа с диска в представление в памяти.
Преобразование кортежей можно ускорить, создав функции, предназначенные для определённой структуры таблицы и количества извлекаемых столбцов.
- Для сокращения затрат JIT компиляция может встраивать (inline) код функций непосредственно в код выражений, вызывающих эти функции.
- LLVM поддерживает оптимизации скомпилированного машинного кода

jit_above_cost

Устанавливает предел стоимости запроса, при превышении которого включается JIT-компиляция

cost > 100_000



jit_inline_above_cost

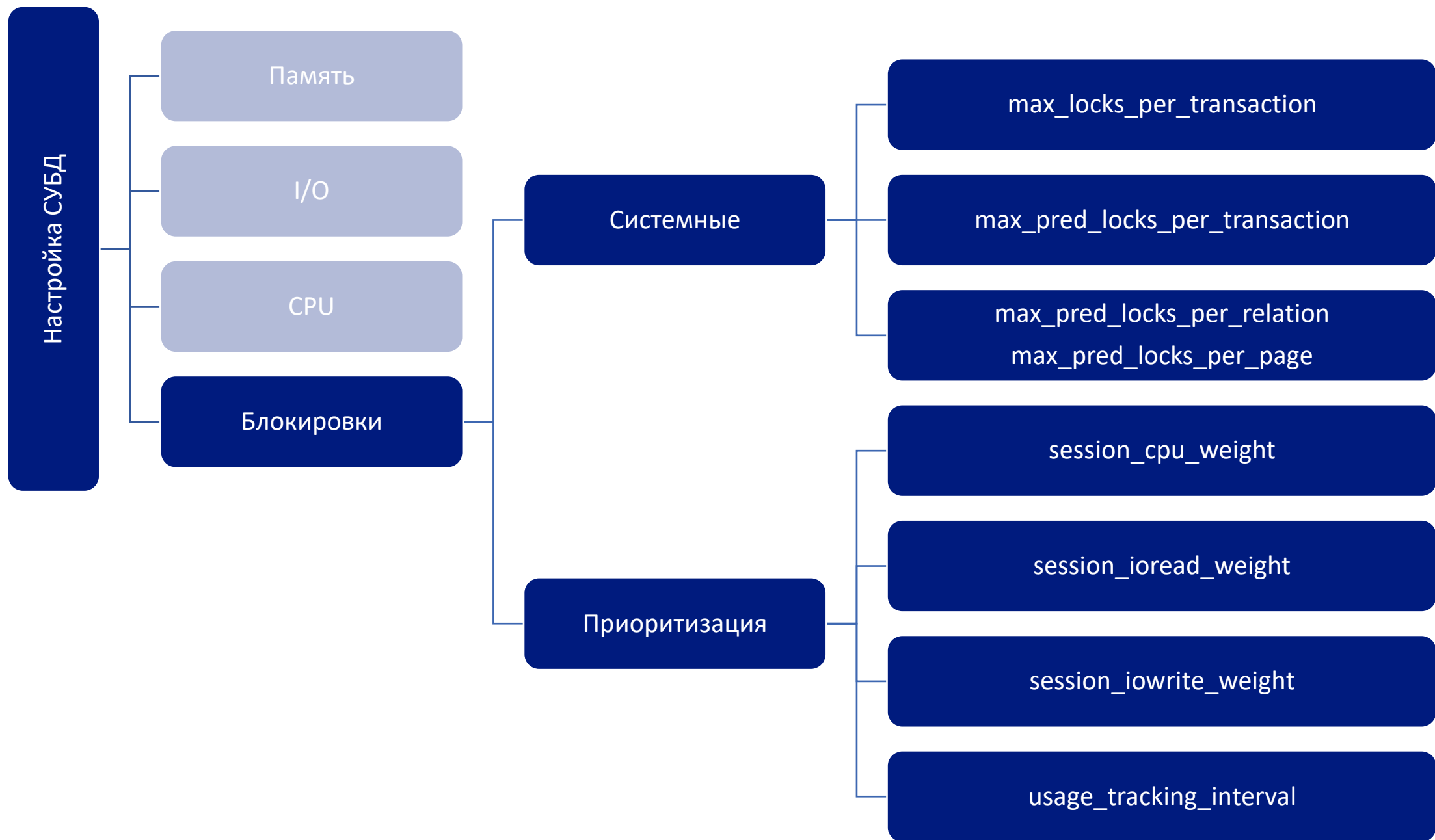
Устанавливает предел стоимости, при превышении которого будет допускаться встраивание функций и операторов в процессе JIT-компиляции.

jit_optimize_above_cost

Устанавливает предел стоимости, при превышении которого в JIT-компилированных программах может применяться дорогостоящая оптимизация.

cost > 500_000

OLTP — отключать
HTAP\OLAP — возможно включать



Пул блокировок

Информация в общей памяти сервера

представление pg_locks:

locktype — тип блокируемого ресурса,

mode — режим блокировки

Ограниченное количество: $\text{max_locks_per_transaction} \times \text{max_connections}$

Реализация уровня изоляции serializable

Типы ресурсов

Relation — отношение

Page — страница

max_pred_locks_per_page



max_pred_locks_per_relation

tuple — версия строки

Режим:

SIRead

- Приоритеты для пользователей и сессий при ограниченном ресурсе
- Методы выделения ресурсов:
 - CPU, чтение, запись
 - Приоритеты 1, 2, 4, 8
 - По умолчанию - 4
- Группы потребителей
 - пользователь может быть членом нескольких групп
 - только одна группа активна во время сеанса
 - группа по умолчанию назначается пользователю при соединении с БД
- Планы и приоритеты можно менять не останавливая PG
- Автоматическое назначение в группу при login
- Автоматическое/ручное переключение планов
- Сеансы работают без ограничений, если ресурсы свободны

WEEKDAY_PLAN

Resource Consumer Group	Resource Plan Directives
Order Entry	CPU = 1 Приоритет по чтению = 4
Shipping	CPU = 2 Приоритет по чтению = 4

NIGHT_PLAN

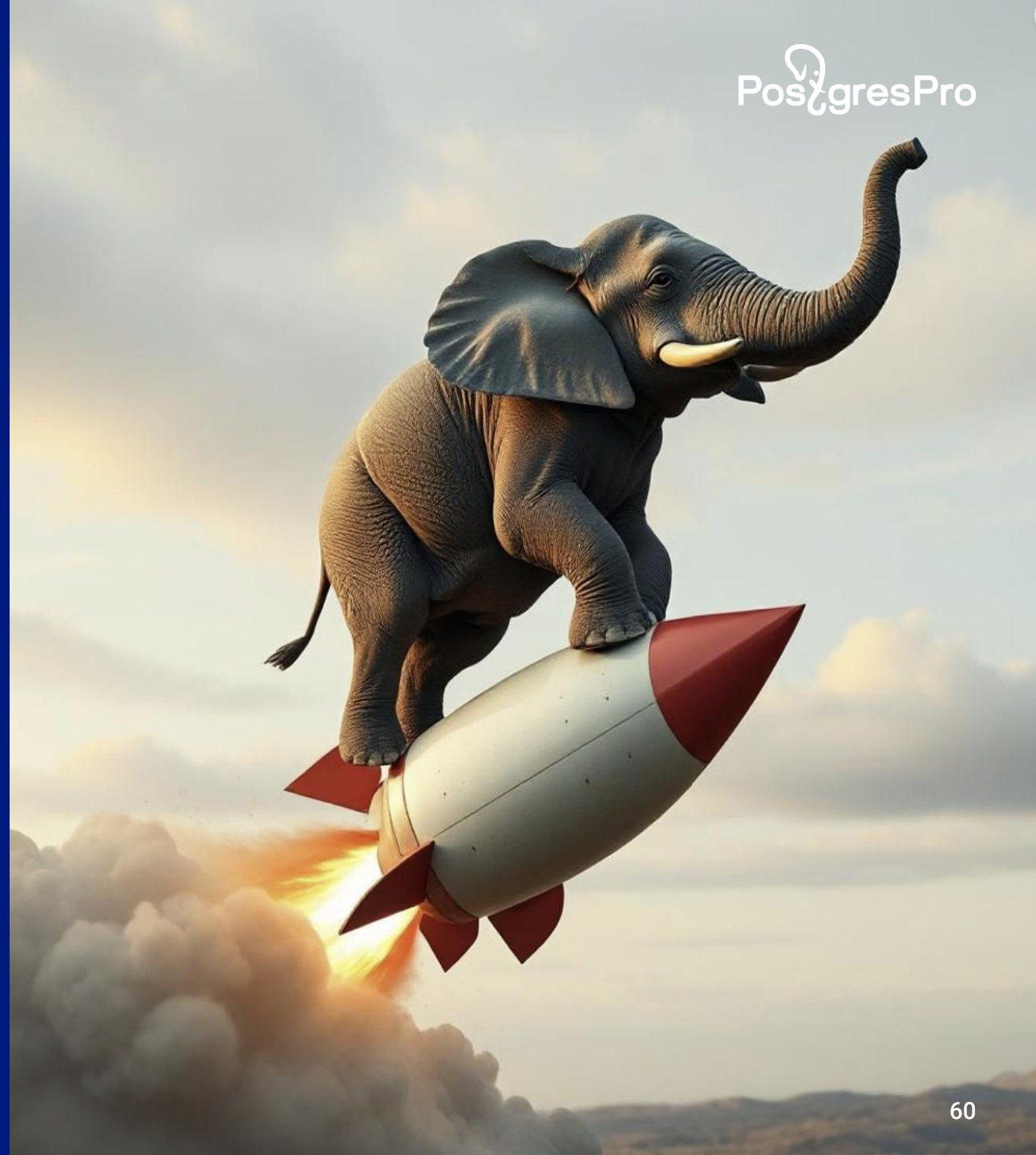
Resource Consumer Group	Resource Plan Directives
Billing	CPU = 4 Приоритет по чтению = 4
Order Entry	CPU 2 Приоритет по записи = 2
Shipping	CPU = 2 Приоритет по чтению = 2

pgpro_tune

База данных на максимум одним нажатием.

Особенно критична становится эта проблема при нехватке времени или необходимости быстро масштабироваться. Ручная настройка может превратиться в серьёзное узкое место при развитии инфраструктуры и запуске новых сервисов.

Команда Postgres Professional реализовала возможность автоматически получить конфигурацию Postgres Pro, учитывающую особенности вашего оборудования и нагрузок, всего за несколько минут, избегая при этом длительных погружений в документацию и многочасовые эксперименты. Знакомьтесь — **pgpro_tune**.



Принцип работы

Решение pgpro_tune — «автотюнер» для Postgres Pro, позволяющий автоматически оптимизировать конфигурацию сервера под оборудование и задачи.

Сбор информации о системе.

pgpro_tune анализирует конфигурацию сервера, определяя количество процессорных ядер, объем оперативной памяти и другие важные параметры оборудования.

Вызов экспертной логики.

pgpro_tune передает собранную информацию о системе набору специализированных shell-скриптов, содержащих экспертные знания и логику настройки, разработанные специалистами Postgres Professional.

Формирование итоговой конфигурации.

Shell-скрипты анализируют входные данные и предлагают оптимальные значения для различных параметров конфигурации. Утилита pgpro_tune агрегирует эти предложения и генерирует итоговый блок настроек.

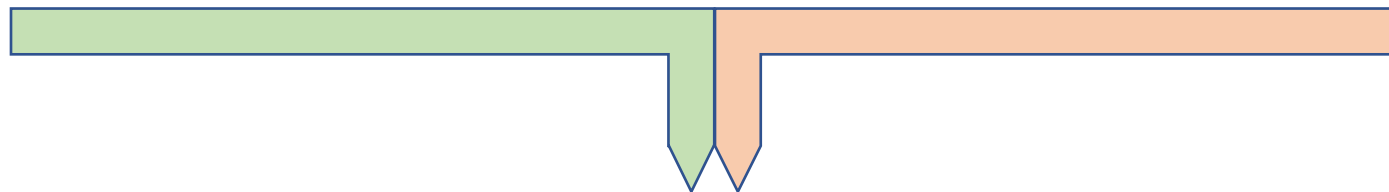
Применение настроек.

Утилита записывает сгенерированные настройки в файл postgresql.conf, снабжая их подробными комментариями и указанием даты и времени автонастройки.

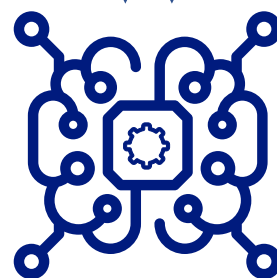


пресеты для прикладного ПО (например, для работы с «1С:Предприятие»). Ориентированы на максимальную производительность в конкретных сценариях использования, учитывая типичные нагрузки и требования прикладного программного обеспечения;

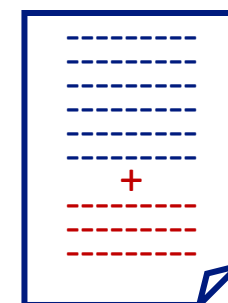
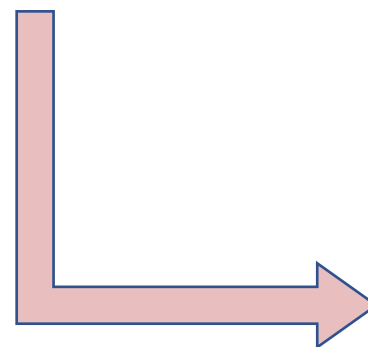
пресеты для задач администрирования. Фокусируются на оптимизации производительности для административных задач (резервное копирование, восстановление или миграция данных)



Сбор информации об аппаратной конфигурации сервера, параметрах ОС, редакции Postgres Pro и тд. (40+ параметров)



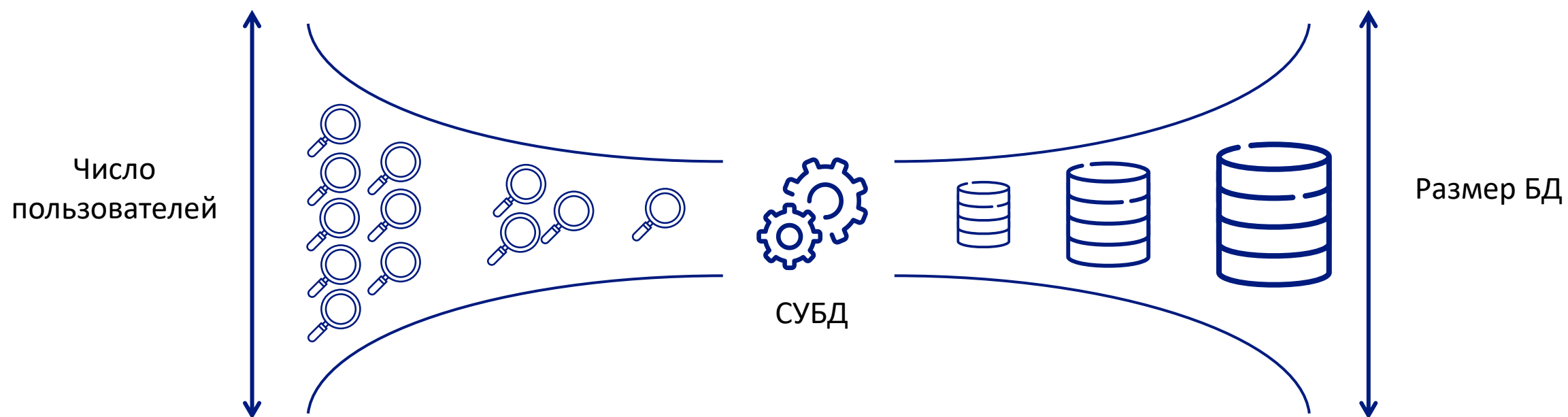
pgpro_tune



postgresql.conf

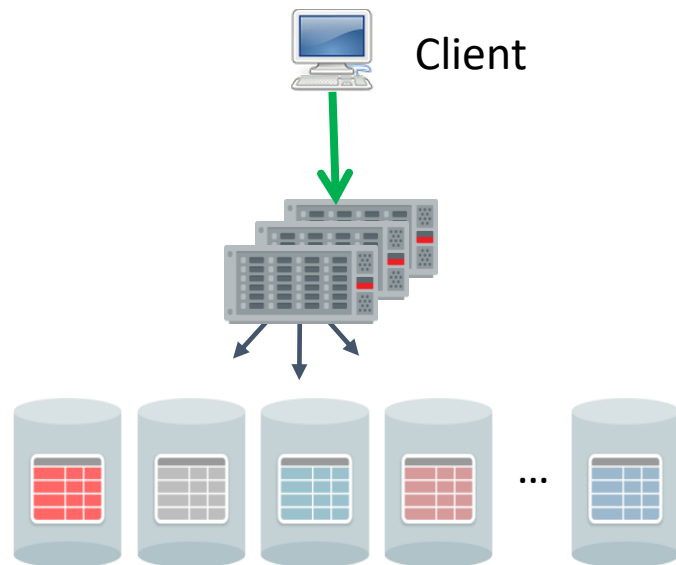


Проблема масштабируемости

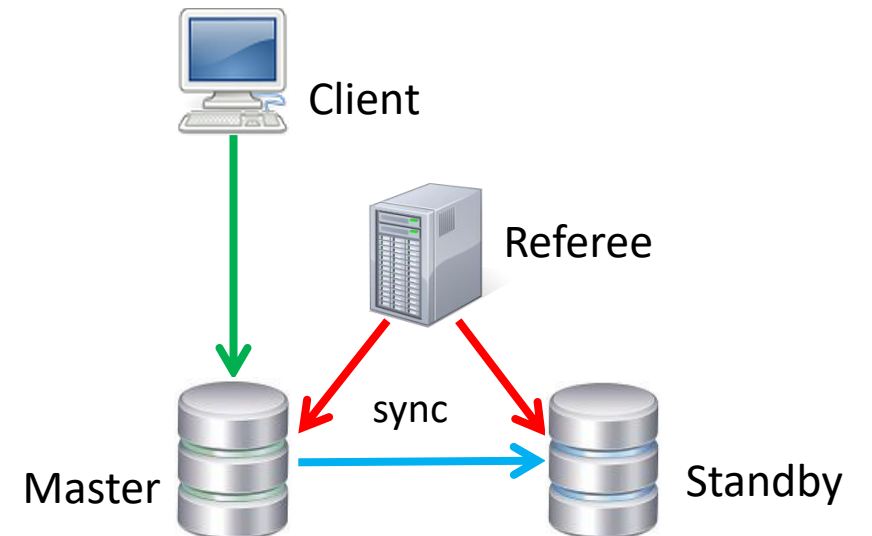


Для высоко нагруженных систем

SHARDMAN

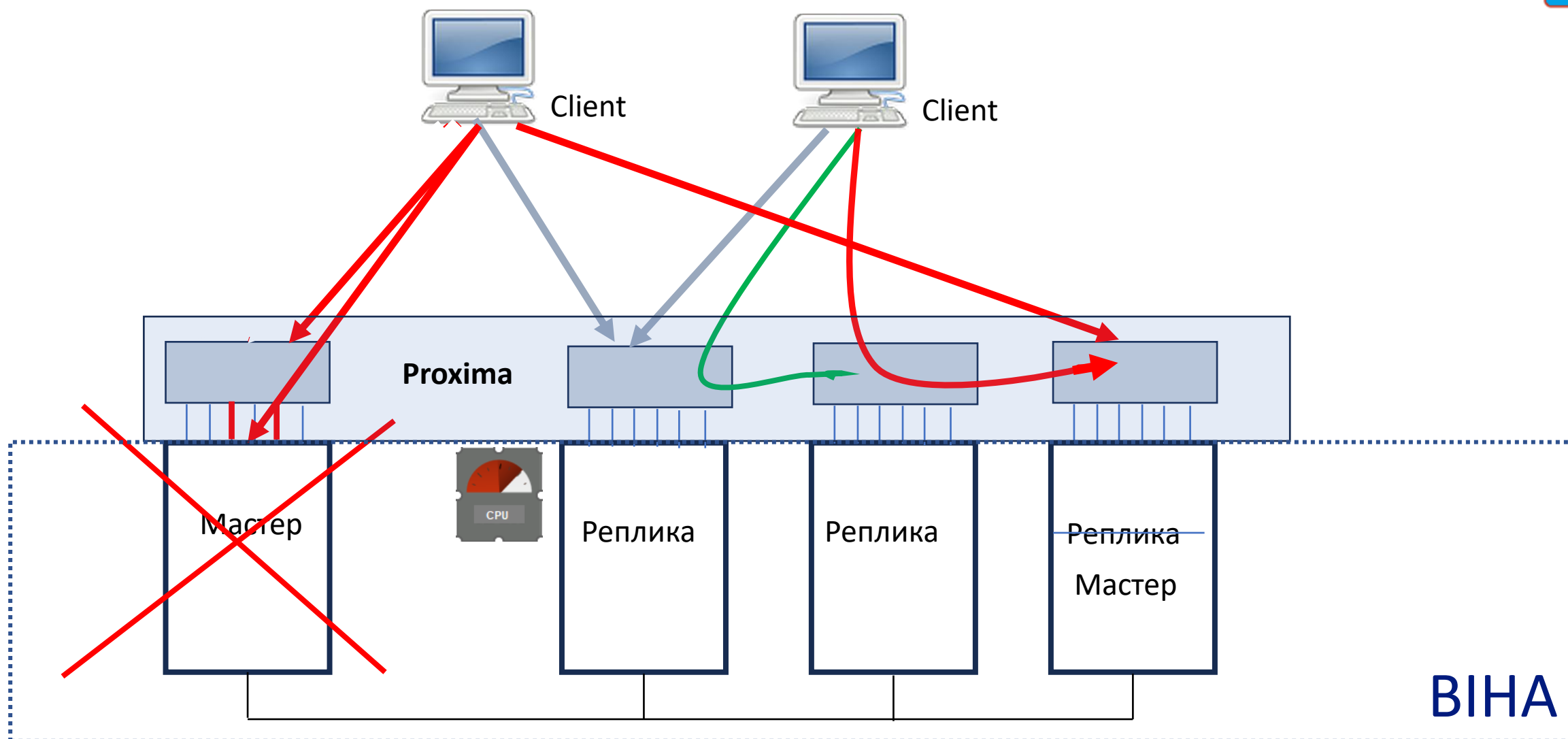


BiHA + proxima load balancer



PROXIMA = Pooler + Proxy + Load Balancer

17.2



Спасибо за внимание!

Q&A