

Where did my memory go?

Frits Hoogland
Yugabyte

PostgreSQL relies on the OS for memory

In most cases 'the OS' will be Linux.

Linux memory usage can be seen in many ways.

A common way is: `/proc/meminfo`

```
[vagrant@pg15 ~]$ cat /proc/meminfo
MemTotal:        1856020 kB
MemFree:         1471340 kB
MemAvailable:    1555788 kB
Buffers:         4204 kB
Cached:          203252 kB
SwapCached:      0 kB
Active:           60820 kB
Inactive:         215696 kB
Active(anon):     784 kB
Inactive(anon):   76888 kB
Active(file):     60036 kB
Inactive(file):   138808 kB
Unevictable:      0 kB
Mlocked:          0 kB
SwapTotal:       2154492 kB
SwapFree:        2154492 kB
Dirty:            32 kB
Writeback:        0 kB
AnonPages:       69064 kB
Mapped:          93892 kB
Shmem:           8612 kB
KReclaimable:    33040 kB
Slab:            68128 kB
SReclaimable:    33040 kB
SUnreclaim:     35088 kB
KernelStack:    2224 kB
PageTables:      5856 kB
NFS_Unstable:    0 kB
Bounce:          0 kB
WritebackTmp:    0 kB
CommitLimit:    3082500 kB
Committed_AS:   324464 kB
VmallocTotal:    34359738367 kB
VmallocUsed:     0 kB
VmallocChunk:    0 kB
Percpu:          992 kB
HardwareCorrupted: 0 kB
AnonHugePages:   14336 kB
ShmemHugePages:  0 kB
ShmemPmdMapped:  0 kB
FileHugePages:   0 kB
FilePmdMapped:   0 kB
HugePages_Total: 0
HugePages_Free:  0
HugePages_Rsvd:  0
HugePages_Surp:  0
Hugepagesize:    2048 kB
Hugetlb:         0 kB
DirectMap4k:    90048 kB
DirectMap2M:    2007040 kB
```

PostgreSQL relies on the OS for memory

Problem: lots of entries.

/proc/meminfo is a motley gathering of memory related statistics.

Multiple views: (in)active, file/anon, memory/swap, hugepages (pages!), virtual, etc.

```
[vagrant@pg15 ~]$ cat /proc/meminfo
MemTotal:      1856020 kB
MemFree:       1471340 kB
MemAvailable:  1555788 kB
Buffers:       4204 kB
Cached:        203252 kB
SwapCached:    0 kB
Active:        60820 kB
Inactive:      215696 kB
Active(anon):  784 kB
Inactive(anon): 76888 kB
Active(file):  60036 kB
Inactive(file): 138808 kB
Unevictable:   0 kB
Mlocked:       0 kB
SwapTotal:     2154492 kB
SwapFree:      2154492 kB
Dirty:         32 kB
Writeback:     0 kB
AnonPages:     69064 kB
Mapped:        93892 kB
Shmem:         8612 kB
KReclaimable:  33040 kB
Slab:          68128 kB
SReclaimable:  33040 kB
SUnreclaim:   35088 kB
KernelStack:  2224 kB
PageTables:    5856 kB
NFS_Unstable:  0 kB
Bounce:        0 kB
WritebackTmp:  0 kB
CommitLimit:  3082500 kB
Committed_AS: 324464 kB
VmallocTotal: 34359738367 kB
VmallocUsed:   0 kB
VmallocChunk:  0 kB
Percpu:        992 kB
HardwareCorrupted: 0 kB
AnonHugePages: 14336 kB
ShmemHugePages: 0 kB
ShmemPmdMapped: 0 kB
FileHugePages: 0 kB
FilePmdMapped: 0 kB
HugePages_Total: 0
HugePages_Free: 0
HugePages_Rsvd: 0
HugePages_Surp: 0
Hugepagesize: 2048 kB
Hugetlb:       0 kB
DirectMap4k:  90048 kB
DirectMap2M:  2007040 kB
```

Linux allows memory oversubscription

Setting: `vm.overcommit_memory`

Default: 0 (heuristic overcommit) -> OOM killer.

Best: 2 (don't overcommit) -> Error on memory allocation (single backend).

Test!

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

Shmem (SysV /dev/zero mmapped)

Dirty

SwapTotal

SwapFree

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

Shmem (SysV /dev/zero mmapped)

Dirty

SwapTotal

SwapFree

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

Shmem (SysV /dev/zero mmapped)

Dirty

SwapTotal

SwapFree

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

Shmem (SysV /dev/zero mmapped)

Dirty

SwapTotal

SwapFree

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

Shmem (SysV /dev/zero mmapped)

Dirty

SwapTotal

SwapFree

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

Shmem (SysV /dev/zero mmapped -- ??)

Dirty

SwapTotal

SwapFree

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

Shmem (SysV /dev/zero mmapped)

Dirty

SwapTotal

SwapFree

Linux - meminfo - Memory hierarchy v1.0

Virtual Memory (not a statistic)

MemTotal

MemFree

Slab

PageTables

KernelStack

SwapCached

AnonPages

Cached

Mapped

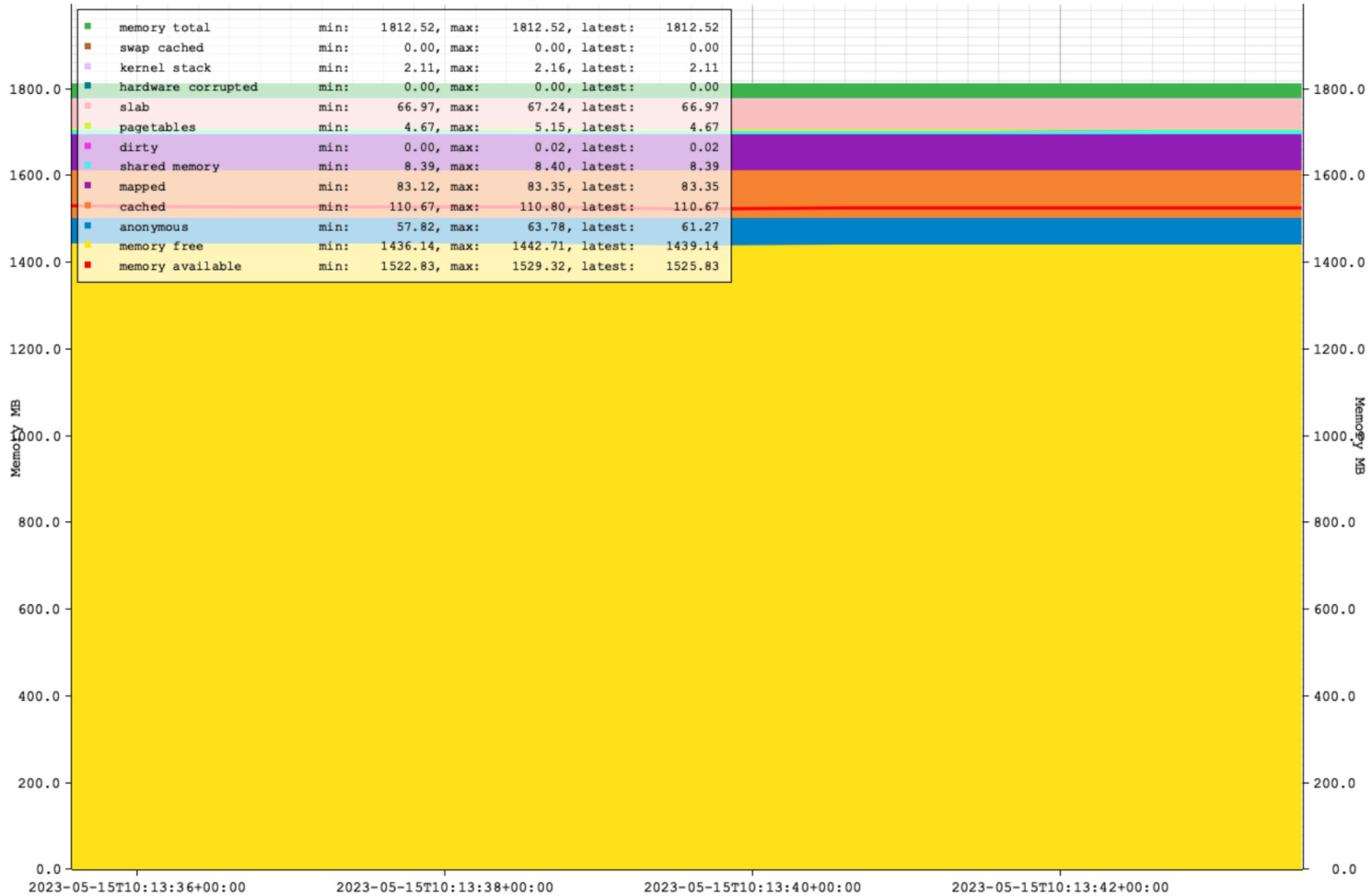
Shmem (SysV /dev/zero mmapped)

Dirty

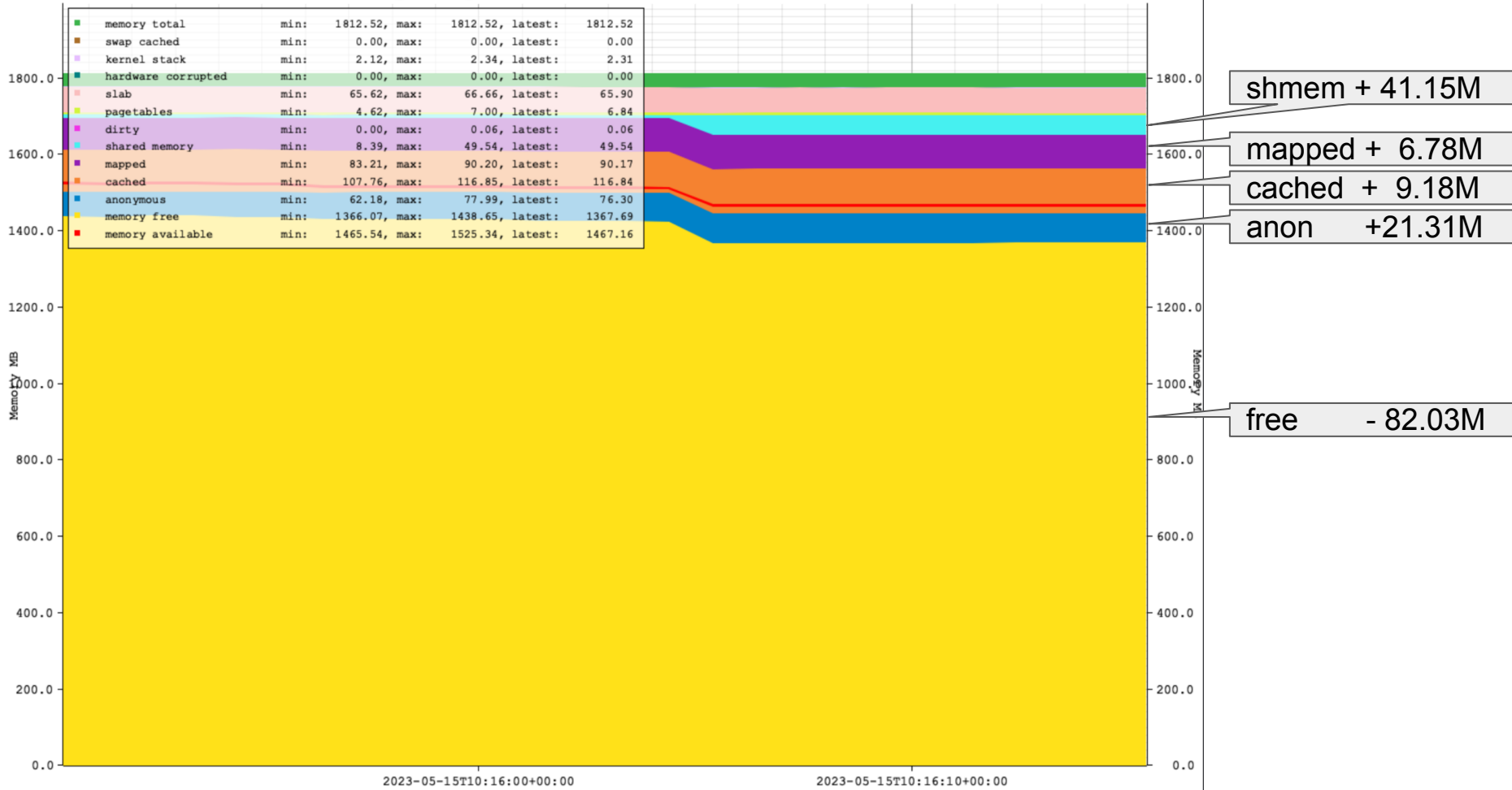
SwapTotal

SwapFree

Memory usage: 192.168.66.100:9300:metrics



Memory usage: 192.168.66.100:9300:metrics



The shared memory segment from 'smaps':

```
7f5187787000-7f51ca8d3000 rw-s 00000000 00:01 33813 /dev/zero (deleted)
Size: 1099056 kB
KernelPageSize: 4 kB
MMUPageSize: 4 kB
Rss: 41012 kB
Pss: 26659 kB
Shared_Clean: 0 kB
Shared_Dirty: 25540 kB
Private_Clean: 0 kB
Private_Dirty: 15472 kB
Referenced: 41012 kB
Anonymous: 0 kB
LazyFree: 0 kB
AnonHugePages: 0 kB
ShmemPmdMapped: 0 kB
FilePmdMapped: 0 kB
Shared_Hugetlb: 0 kB
Private_Hugetlb: 0 kB
Swap: 0 kB
SwapPss: 0 kB
Locked: 0 kB
THPeligible: 0
VmFlags: rd wr sh mr mw me ms sd
```

Let's test!

```
-- precreated table:
create table t( id int primary key, f1 text);
insert into t select id, repeat('x',16384) from generate_series(1,4500000) id;
```

```
-- table size
```

```
postgres=# \dt+ t
```

List of relations

Schema	Name	Type	Owner	Persistence	Access method	Size	Description
public	t	table	postgres	permanent	heap	1034 MB	

```
-- index size
```

```
postgres=# \di+ t_pkey
```

List of relations

Schema	Name	Type	Owner	Table	Persistence	Access method	Size	Description
public	t_pkey	index	postgres	t	permanent	btree	96 MB	

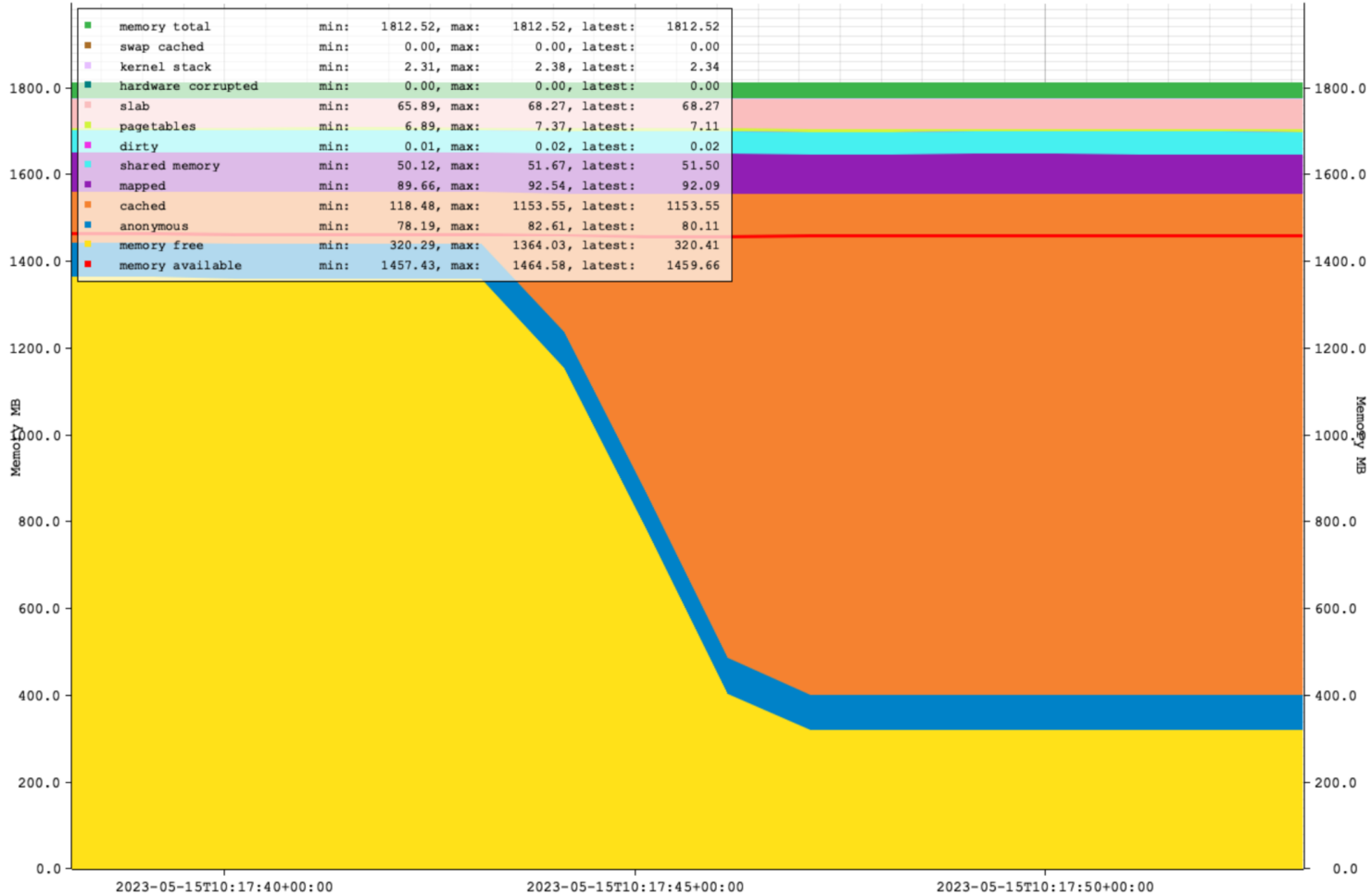
1. Force seq scan

```
explain analyze select count(f1) from t;
```

QUERY PLAN

```
-----  
Finalize Aggregate (cost=156790.72..156790.73 rows=1 width=8) (actual time=2395.054..2403.280 rows=1 loops=1)  
  -> Gather (cost=156790.50..156790.71 rows=2 width=8) (actual time=2394.649..2403.272 rows=3 loops=1)  
      Workers Planned: 2  
      Workers Launched: 2  
      -> Partial Aggregate (cost=155790.50..155790.51 rows=1 width=8) (actual time=2367.083..2367.083 rows=1 loops=1)  
          -> Parallel Seq Scan on t (cost=0.00..151103.00 rows=1875000 width=200) (actual time=0.614..0.614 rows=1875000 loops=1)  
Planning Time: 7.306 ms  
Execution Time: 2403.614 ms  
(8 rows)
```

Memory usage: 192.168.66.100:9300:metrics



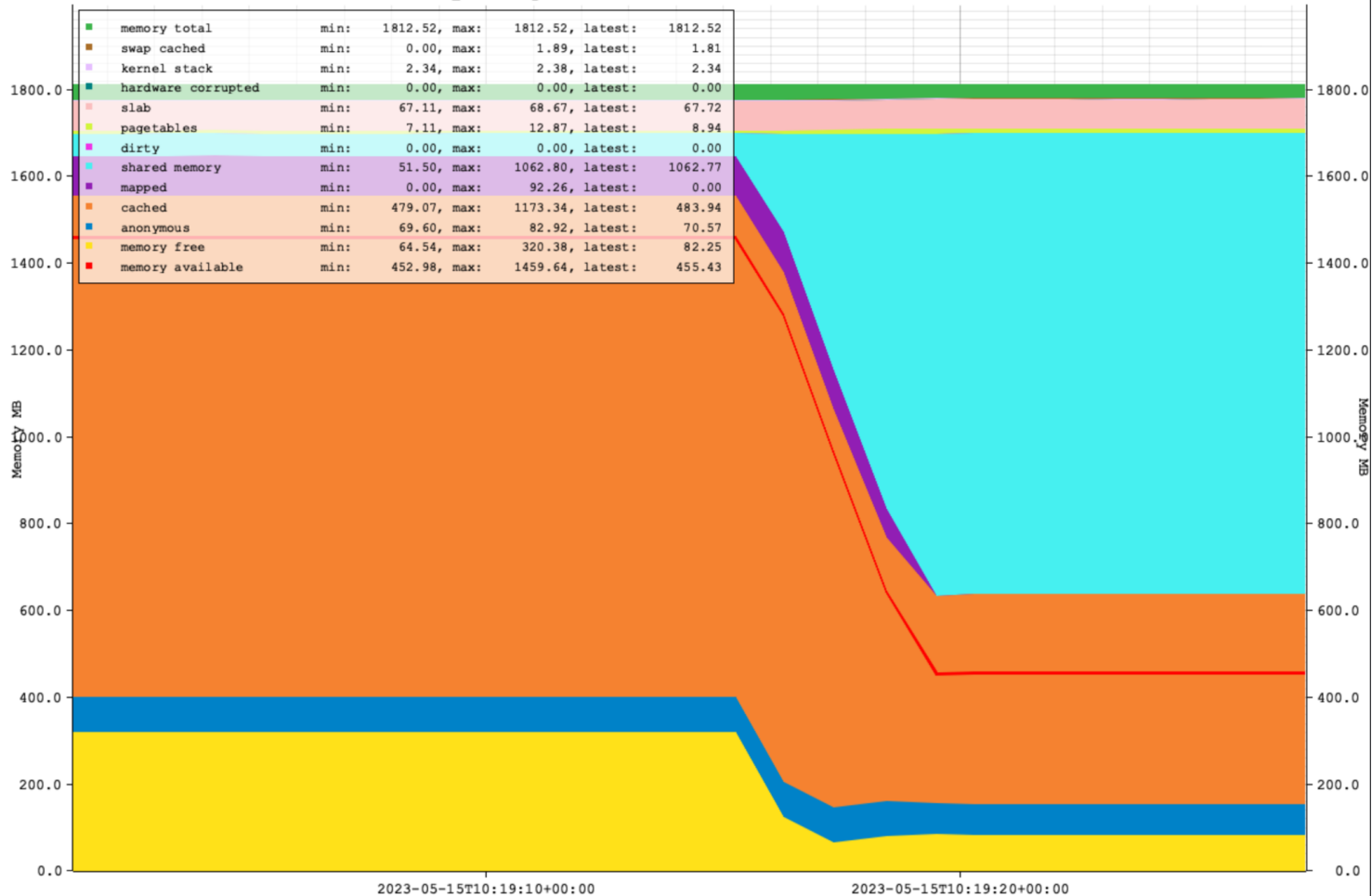
2. Force index scan

```
set enable_seqscan=off;  
explain analyze select count(f1) from t where id > 0;
```

QUERY PLAN

```
Index Scan using t_pkey on t  (cost=0.43..260474.43 rows=4500000 width=204) (actual time=0.036..4414.006  
  Index Cond: (id > 0)  
Planning Time: 6.661 ms  
Execution Time: 4716.215 ms  
(4 rows)
```

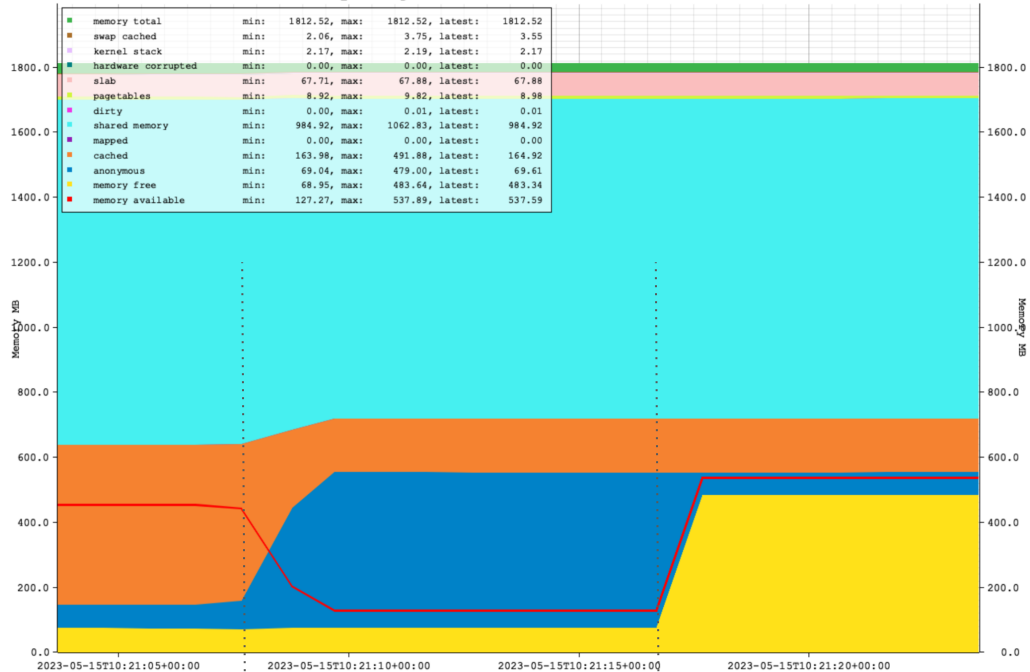
Memory usage: 192.168.66.100:9300:metrics



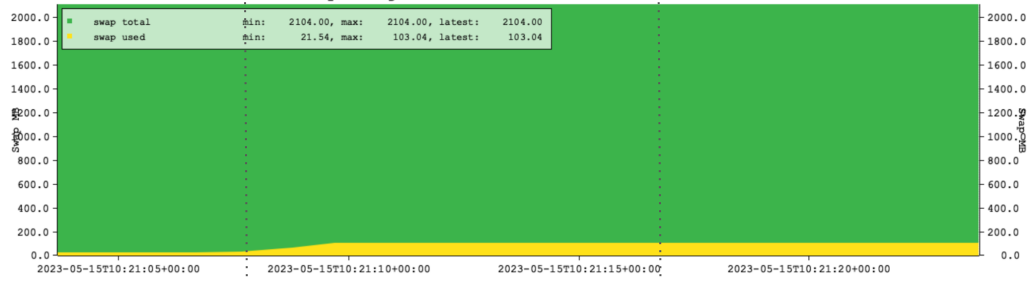
2. Force backend allocation

```
set my.count to 409600;
set my.size to 1020;
do
$$
  declare
    text_array text[];
    counter    int:= current_setting('my.count',true);
    size       int:= current_setting('my.size',true);
    size_read  varchar;
  begin
    raise info 'Pid: %', pg_backend_pid();
    raise info 'Array element size: %, count: %', size, counter;
    for count in 1..counter loop
      text_array[count]:=repeat('x',size);
    end loop;
    raise info 'done!';
    select pg_size_pretty(sum(total_bytes)) into size_read from pg_backend_memory_contexts;
    raise info 'size now: %', size_read;
    perform pg_sleep(60);
  end
$$;
```

Memory usage: 192.168.66.100:9300:metrics



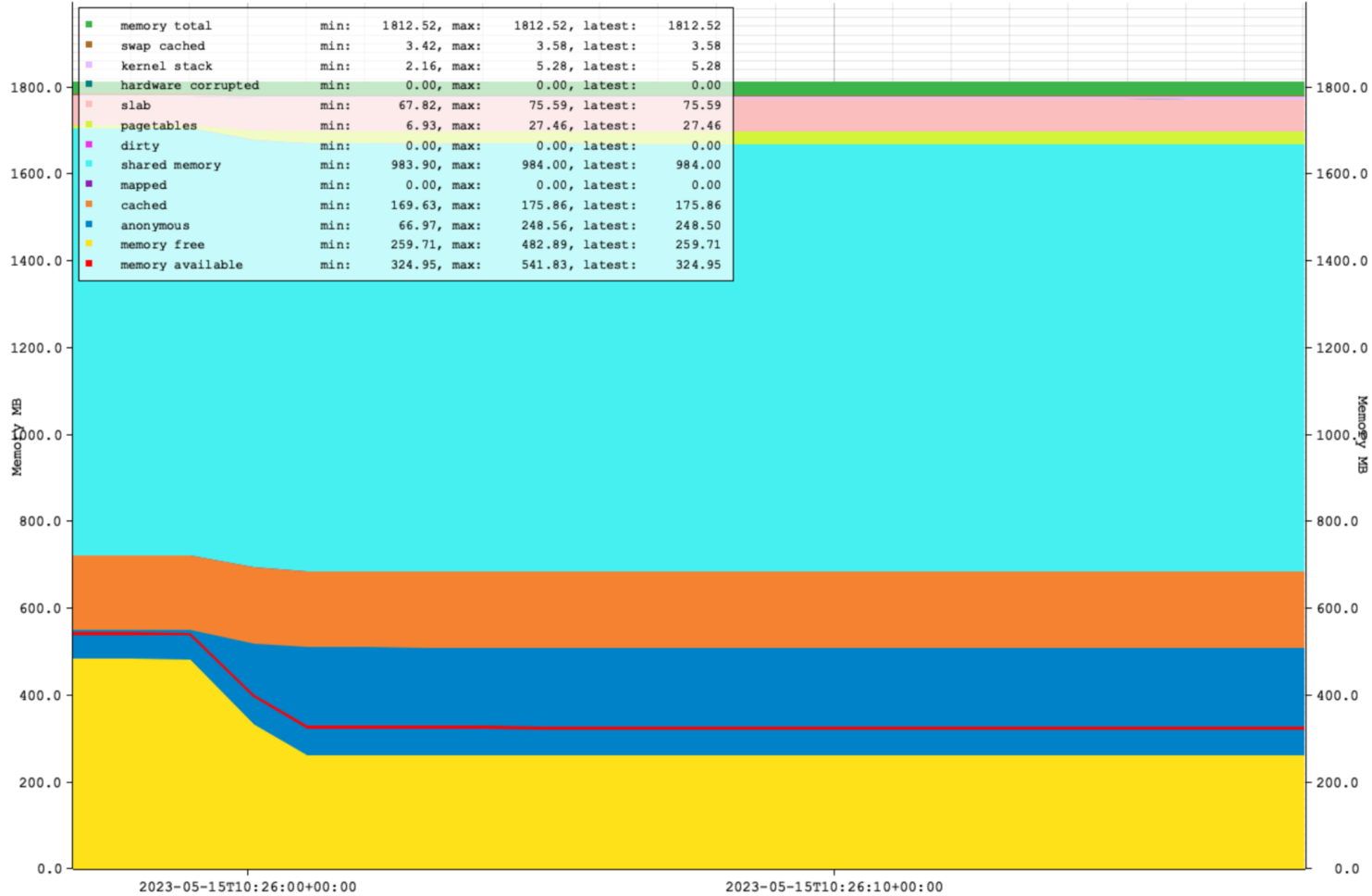
Swap usage: 192.168.66.100:9300:metrics



3.a Lots of backends -- idle

```
for T in $(seq 1 100); do
    psql -c "select pg_sleep(60)" &
done
```

Memory usage: 192.168.66.100:9300:metrics



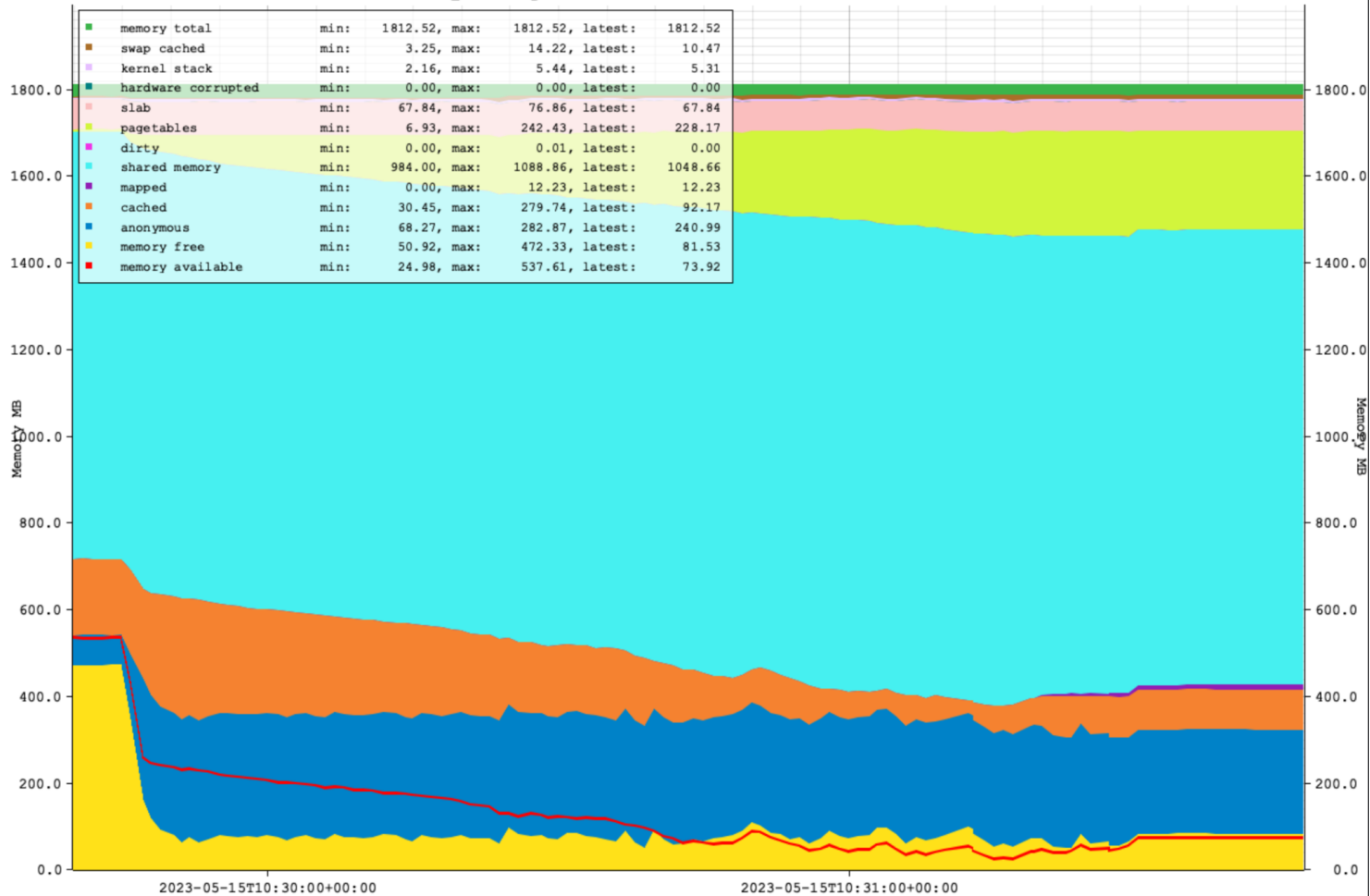
3.a Lots of backends -- idle

```
postgres=# select pg_size_pretty(sum(total_bytes)) from pg_backend_memory_contexts;
pg_size_pretty
-----
1862 kB
```

3.b Lots of backends -- just doing simple things

```
for T in $(seq 1 100); do
    psql -c "set enable_seqscan=off; select count(f1) from t where id > 0; select pg_sleep(60)" &
done
```

Memory usage: 192.168.66.100:9300:metrics



4. Large backend allocations

Here is an oddity I found that could be very confusing.

In order to show this:

- I restart the server to decrease allocated shared memory.
- This allows anonymous allocations to be the majority of memory allocated.
- I run the PLpgsql array allocation procedure with different settings:

4. Large backend allocations

```
select pg_size_pretty(sum(total_bytes)) from pg_backend_memory_contexts;
set my.count to 3000;
set my.size to 166666;
do
$$
  declare
    text_array text[];
    counter    int:= current_setting('my.count',true);
    size       int:= current_setting('my.size',true);
    size_read  varchar;
  begin
    raise info 'Pid: %', pg_backend_pid();
    raise info 'Array element size: %, count: %', size, counter;
    for count in 1..counter loop
      text_array[count]:=repeat('x',size);
    end loop;
    raise info 'done!';
    select pg_size_pretty(sum(total_bytes)) into size_read from pg_backend_memory_contexts;
    raise info 'size now: %', size_read;
    perform pg_sleep(5);
  end
$$;
select pg_size_pretty(sum(total_bytes)) from pg_backend_memory_contexts;
```

4. Large backend allocations

```
postgres=# \i alloc_x.sql  
pg_size_pretty
```

```
-----
```

```
1755 kB  
(1 row)
```

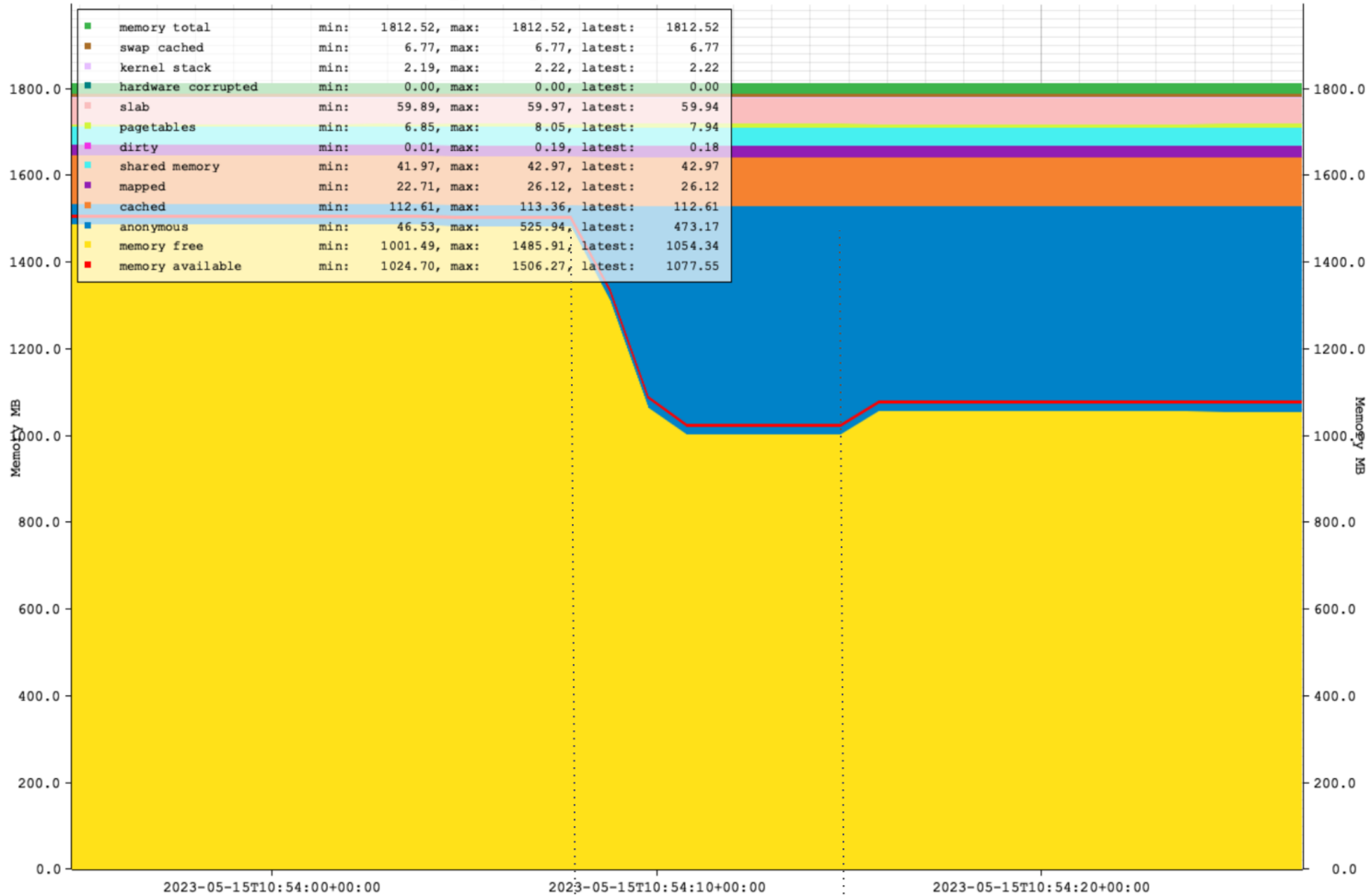
```
SET  
SET  
psql:alloc_x.sql:56: INFO:  Pid: 8646  
psql:alloc_x.sql:56: INFO:  Array element size: 166666, count: 3000  
psql:alloc_x.sql:56: INFO:  done!  
psql:alloc_x.sql:56: INFO:  size now: 479 MB
```

```
DO  
pg_size_pretty
```

```
-----
```

```
1774 kB  
(1 row)
```

Memory usage: 192.168.66.100:9300:metrics



4. Large backend allocations – short explanation

Anonymous allocations in PostgreSQL are performed by `palloc()`.
Palloc is a PostgreSQL abstraction, and calls the OS-level `malloc()` call.

`malloc()` is not a system call. It is a **library** to manage allocations in an address space.

So it's a memory management software layer between PostgreSQL and the OS.

It keeps allocated chunks around/cached 'in case you need it'.

- Small enough allocations are hardly visible/influencing.
- Very large allocations will get freed.
- Medium to large allocations are cached and can be significant.

Malloc resources

Malloc internals: <https://sourceware.org/glibc/wiki/MallocInternals>

Azeria heap exploitation part 1 (intro): <https://azeria-labs.com/heap-exploitation-part-1-understanding-the-glibc-heap-im>

Azeria heap exploitation part 2 (free): <https://azeria-labs.com/heap-exploitation-part-2-glibc-heap-free-bins/>

Conclusion

- Memory statistics in linux require interpretation and calculation.
 - I am not sure mine is 100% correct, but it seems to give a reasonable impression.
 - Any updates or corrections are welcome!
- Memory allocations in linux are lazy allocated.
 - And therefore PostgreSQL allocations are lazy allocated too.
 - Notably the buffer cache; anything else in the backend is even more dynamic.
- For architects and testers:
 - This means you need to build quite sophisticated test cases with real life data to be realistic.
 - Too simple test cases mean you actually/probably test a different scenario.

Conclusion

- Free memory != usable memory.
 - Look at available memory.
 - Too much free memory consistently reported likely means not using all memory.
- Main common PostgreSQL allocations:
 - Shared memory (== mapped memory == cached memory): buffer cache.
 - Anonymous memory: backend memory allocation.
- PostgreSQL is exhaustively documented to be needing OS cache.
 - Actual cache == cached - shared memory - mapped memory - dirty pages
 - Do you leave memory for OS caching?
 - If not: this will make query latency fluctuate a lot for seemingly identical situations.

Conclusion

- Large seq scans will only use 256k of buffers as ring.
 - Large means > 25% of set buffer cache size.
- Backends require a low amount of memory.
 - But need to build a pagetables to access buffer cache.
 - Any actual work requiring backend memory allocation comes on top.
 - Therefore: be very reserved about the number of backends.
- Malloc can play memory tricks.
 - For small allocations this is probably as good as invisible.
 - For 'medium' allocations this might play up.
 - I am not sure how often this plays up, but it's hard to spot if you're unaware.

Conclusion

- dsar (which created the memory graphs):
 - <https://github.com/fritshoogland-yugabyte/dsar>
 - Any requests are welcome!