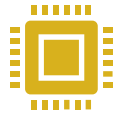# About me

**EM for PG OSS contributors/committers team.**

**Started with porting Linux to new ;-) MIPS, SH4 boards and device drivers !!**

**~22 years industry experience**

Naturesoft (Embedded)

HP (HPC, Parallel filesystems)

Storsimple (Hybrid cloud storage)

Microsoft (Storsimple, Kubernetes, Flex PG service, OSS PG)
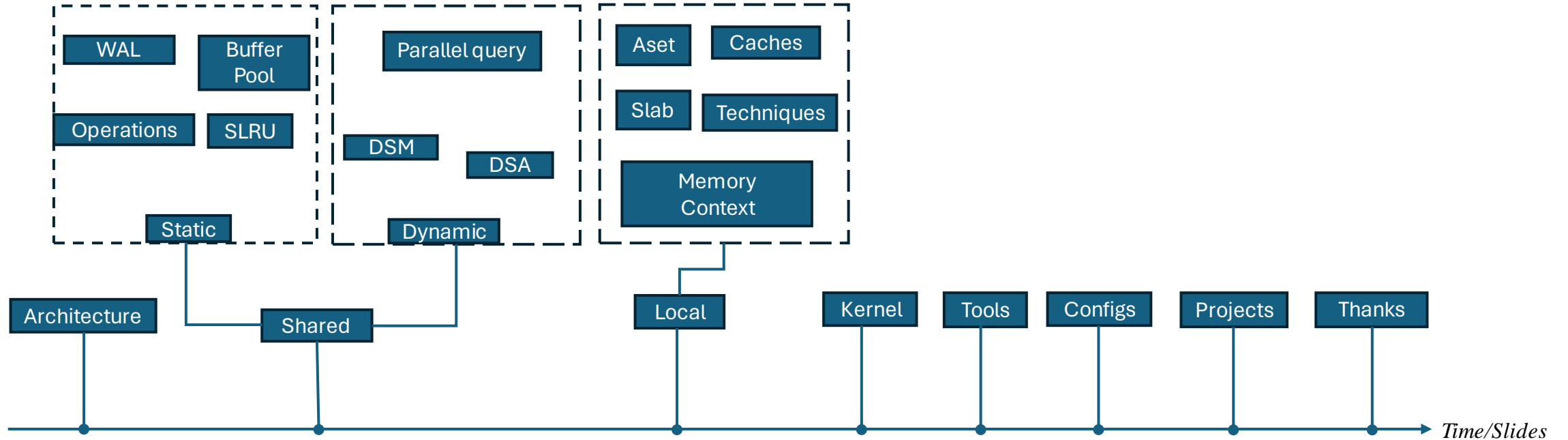
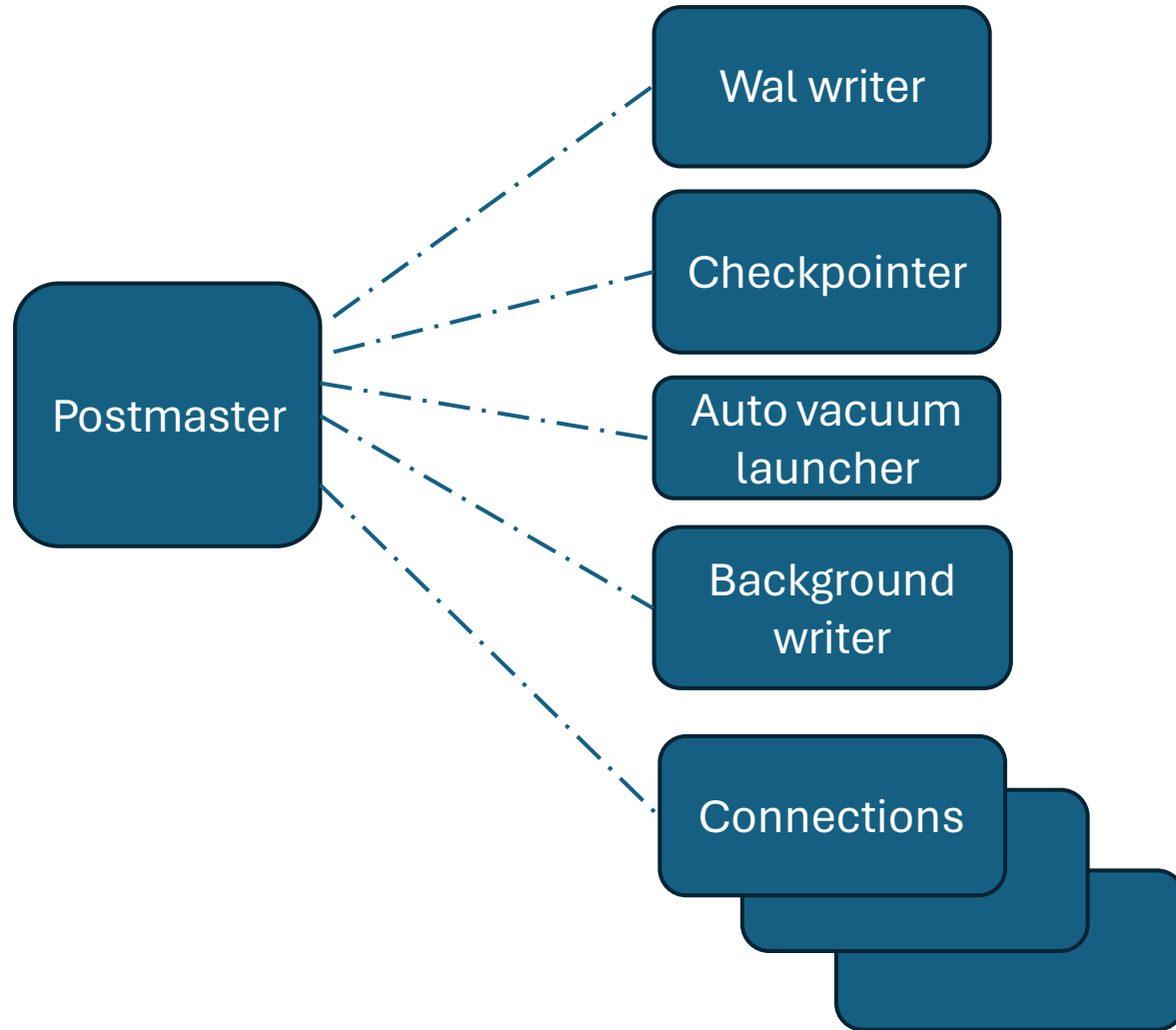**One of the most exciting part of my day – interacting with top PostgreSQL folks.**

**Have come to share the lessons learned on PostgreSQL memory management.**

# What will we discuss?

# Process based architecture



Postmaster

- Wal writer
- Checkpointer
- Auto vacuum launcher
- Background writer
- Connections

Currently, POSTGRES runs as one process for each active user. This was done as an expedient to get a system operational as quickly as possible. We plan on converting POSTGRES to use lightweight processes available in the operating systems we are using

- *"The implementation of POSTGRES". M. Stonebraker, L. A. Rowe, and M. Hirohama. Transactions on Knowledge and Data Engineering 2(1). IEEE. March 1990.*

# Memory types

- Shared
- Local
- Kernel

# Shared memory

- Static
- Dynamic

# Static memory

- Create
- Locking
- Examples:
  - Buffer pool
  - SLRU
  - WAL buffer

# Creation

- SysV memory
- Mmap

# Lightweight Locking – slow path

P1 → *LWLockAcquire()*

*LWLockQueueSelf ()*

*LWLockAttemptLock ()*

*PGSemaphoreLock ()*

*Queue*

*Atomic memory operations*

*Wait*

**Shared memory**

Lock wait Q

Semaphore

Lock

P2

- Short lived situations in contrast with heavy weigh locks
- DB literature it's generally called latch
- Shared or exclusive
- No deadlock detection
- In future could be built on top of futex ?

# Buffer pool

Clocksweep

LastFree

Buffer Strategy Control

Next Victim

FirstFree

Buffer Access Strategy

Clocksweep

Buffer Hash Table

Buffer descriptors

Buffer blocks

Most frequently & recently
BufferAccessStrategy useful for Scan resistance
- Sequential
- Vacuum

CLOCK algorithm 1960s Multics ?
Future? CAR built-in scan resistance

ReadBufferExtended
Usage count goes up to 5

# Bufferpool contents

```
postgres=# select * from pg_buffercache where relfilenode = 24698;
 bufferid | relfilenode | reltablespace | reldatabase | relforknumber | relblocknumber | isdirty | usagecount | pinning_backends
----------+-------------+---------------+-------------+---------------+----------------+---------+------------+------------------
        1 |       24698 |          1663 |           5 |             0 |         108492 | t       |          5 |                0
        3 |       24698 |          1663 |           5 |             0 |         108493 | t       |          5 |                0
        4 |       24698 |          1663 |           5 |             0 |         108494 | t       |          5 |                0
        5 |       24698 |          1663 |           5 |             0 |         108495 | t       |          5 |                0
        6 |       24698 |          1663 |           5 |             0 |         108496 | t       |          5 |                0
        7 |       24698 |          1663 |           5 |             0 |         108497 | t       |          5 |                0
        9 |       24698 |          1663 |           5 |             0 |         108498 | t       |          5 |                0
       10 |       24698 |          1663 |           5 |             0 |         108499 | t       |          5 |                0
       11 |       24698 |          1663 |           5 |             0 |         108500 | t       |          5 |                0
       12 |       24698 |          1663 |           5 |             0 |         108501 | t       |          5 |                0
       13 |       24698 |          1663 |           5 |             0 |         108502 | t       |          5 |                0
       15 |       24698 |          1663 |           5 |             0 |         108503 | t       |          5 |                0
       16 |       24698 |          1663 |           5 |             0 |         108504 | t       |          5 |                0
       17 |       24698 |          1663 |           5 |             0 |         108505 | t       |          5 |                0
       18 |       24698 |          1663 |           5 |             0 |         108506 | t       |          5 |                0
       19 |       24698 |          1663 |           5 |             0 |         108507 | t       |          5 |                0
       21 |       24698 |          1663 |           5 |             0 |         108508 | t       |          5 |                0
       22 |       24698 |          1663 |           5 |             0 |         108509 | t       |          5 |                0
       23 |       24698 |          1663 |           5 |             0 |         108510 | t       |          5 |                0
```

# Other shared memory

- WAL buffer
- SLRU

```
postgres=# select name from pg_stat_slru;
      name
--------------------
 commit_timestamp
 multixact_member
 multixact_offset
 notify
 serializable
 subtransaction
 transaction
 other
(8 rows)
```

# DSM – Dynamic shared memory

```
postgres=# set debug_parallel_query = on;
SET
postgres=# select 2024;
 ?column?
----------
     2024
(1 row)

postgres=# select 2024;
 ?column?
----------
     2024
(1 row)

postgres=#
```

```
$strace -f -p $CONNECTED_PID  2>&1 | egrep "mmap|munmap|shm_open|shm"

openat(AT_FDCWD, "/dev/shm/PostgreSQL.3184094726", O_RDWR|O_CREAT|O_EXCL|O_NOFOLLOW|O_CLOEXEC, 0600) = 7
mmap(NULL, 94784, PROT_READ|PROT_WRITE, MAP_SHARED, 7, 0) = 0x7f5caf47d000
munmap(0x7f5caf47d000, 94784)              = 0
unlink("/dev/shm/PostgreSQL.3184094726") = 0
```

# DSA – Dynamic shared memory areas

- Memory allocator built on top of DSM
- Primarily started for parallel hash join.

# Process local memory

- Interface - MemoryContext
- MemoryContext{Alloc, Realloc, Reset, Delete}
- Implementations
  - Allocation Set - standard
  - Slab (large equally sized objects) – logical replication
- CurrentMemoryContext always point to current

# Hierarchy and error handling

- Hierarchical context
  - Top Memory Context
  - Cache Memory context
- Parent and Child free interlinked
- Exception
  - Set jump
  - Memory freed on error
- Child parent
  - Set as parent after allocations

```
             context_name          | level
-----------------------------------+-------
TopMemoryContext                   |   1
  CacheMemoryContext               |   2
  ErrorContext                     |   2
  GUCMemoryContext                 |   2
  LOCALLOCK hash                   |   2
  MdSmgr                           |   2
  MessageContext                   |   2
  Operator class cache             |   2
  Operator lookup cache            |   2
  PgStat Pending                   |   2
  PgStat Shared Ref                |   2
  PgStat Shared Ref Hash           |   2
  Portal hash                      |   2
  PrivateRefCount                  |   2
  Record information cache         |   2
  Relcache by OID                  |   2
  RowDescriptionContext            |   2
  Timezones                        |   2
  TopPortalContext                 |   2
  TopTransactionContext            |   2
  TransactionAbortContext          |   2
  Type information cache           |   2
  WAL record construction          |   2
  search_path processing cache     |   2
  smgr relation table              |   2
    GUC hash table                 |   3
    PortalContext                  |   3
    index info                     |   3
    relation rules                 |   3
      ExecutorState                |   4
        ExprContext                |   5
        Table function arguments   |   5
        TupleSort main             |   5
        printtup                   |   5
          TupleSort sort           |   6
            Caller tuples          |   7
```

# Kernel memory consumption

- Typical process consumption – page table, stack etc.
- Double buffering in page cache

```
postgres=# SELECT pg_total_relation_size('my_large_table');
 pg_total_relation_size
------------------------
             1074864128
(1 row)

postgres=# select pg_filenode_relation(0, 24698);
 pg_filenode_relation
----------------------
 my_large_table
(1 row)

postgres=# select pg_filenode_relation(0, 24705);
 pg_filenode_relation
----------------------
 my_large_table_pkey
(1 row)
```

```
$fincore   -b 24705 24698
        RES   PAGES       SIZE FILE
 183296000   44750  183296000 24705
 891289600  217600  891289600 24698
$
```

# Configuration

- shared_buffers
- work_mem
    - Executor nodes of the query
    - Parallel workers
    - User sessions
    - Partitions
- huge_pages – on, off, try
- Overcommit settings
    sysctl -w vm.overcommit_memory=2

# Related views

- pg_shmem_allocations
- pg_backend_memory_contexts

```
postgres=# select * from pg_shmem_allocations;
               name                 |    off    |   size    | allocated_size
------------------------------------+-----------+-----------+----------------
 Buffer Descriptors                 |   5730944 |   1048576 |        1048576
 Backend SSL Status Buffer          | 146802560 |     41472 |          41472
 Async Queue Control                | 147476864 |      3952 |           3968
 Wal Sender Ctl                     | 147470080 |      1144 |           1152
 AutoVacuum Data                    | 147461248 |      5328 |           5376
 commit_timestamp                   |   4793472 |    267424 |         267520
 multixact_member                   |   5462400 |    267424 |         267520
 multixact_offset                   |   5328640 |    133760 |         133760
 subtransaction                     |   5061120 |    267424 |         267520
 notify                             | 147480832 |    133760 |         133760
 Shared Memory Stats                | 147614592 |    279992 |         280064
 serializable                       | 146176896 |    267424 |         267520
 PROCLOCK hash                      | 143215104 |      2896 |           2944
 FinishedSerializableTransactions   | 146176768 |        16 |            128
 XLOG Ctl                           |     54912 |   4208192 |        4208256
 Shared MultiXact State             |   5729920 |      1024 |           1024
 Proc Header                        | 146444544 |       136 |            256
 Archiver Data                      | 147473664 |         8 |            128
 XLOG Recovery Ctl                  |   4263680 |       104 |            128
 Backend Client Host Name Buffer    | 146663296 |      8192 |           8192
 ReplicationSlot Ctl                | 147466624 |      2720 |           2816
 KnownAssignedXids                  | 146560128 |     31720 |          31744
 Prepared Transaction Table         | 146844032 |        16 |            128
 BTree Vacuum State                 | 147474560 |      1476 |           1536
 Checkpoint BufferIds               | 141263488 |    327680 |         327680
 Wal Receiver Ctl                   | 147471232 |      2264 |           2304
 PREDICATELOCKTARGET hash           | 143914496 |      2896 |           2944
 Backend Status Array               | 146599808 |     55296 |          55296
 KnownAssignedXidsValid             | 146591872 |      7930 |           7936
 Slot Sync Data                     | 147474432 |        24 |            128
 WaitEventExtensionCounterData      | 147894656 |         8 |            128
 DSM Registry Data                  |     54656 |        16 |            128
 WaitEventExtension hash by name    | 147901184 |      2896 |           2944
 Shared Buffer Lookup Table         | 141591168 |      2896 |           2944
 CommitTs shared                    |   5060992 |        32 |            128
 Backend Application Name Buffer    | 146655104 |      8192 |           8192
 ProcSignal                         | 146925440 |     11272 |          11392
 Logical Replication Launcher Data  | 147473792 |       528 |            640
 Buffer Blocks                      |   6779520 | 134221824 |      134221824
 Buffer IO Condition Variables      | 141001344 |    262144 |         262144
 Proc Array                         | 146559488 |       524 |            640
 PMSignalState                      | 146924416 |      1016 |           1024
 PREDICATELOCK hash                 | 144359808 |      2896 |           2944
 PredXactList                       | 145620224 |        88 |            128
 Fast Path Strong Relation Lock Data| 143910272 |      4100 |           4224
 Wal Summarizer Ctl                 | 147473536 |        48 |            128
 transaction                        |   4263808 |    529568 |         529664
 RWConflictPool                     | 145883776 |        24 |            128
 WaitEventExtension hash by id      | 147894784 |      2896 |           2944
 TransamVariables                   |     54624 |        72 |            128
 XLogPrefetchStats                  |   4263552 |        72 |            128
 Buffer Strategy Status             | 142519808 |        28 |            128
 SerialControlData                  | 146444416 |        12 |            128
 shmInvalBuffer                     | 146856192 |     68128 |          68224
 Sync Scan Locations List           | 147476096 |       656 |            768
```

```
postgres=# select * from pg_backend_memory_contexts limit 20;
         name          |    ident    |        parent        | level | total_bytes | total_nblocks | free_bytes | free_chunks | used_bytes
-----------------------+-------------+----------------------+-------+-------------+---------------+------------+-------------+------------
 TopMemoryContext      |             |                      |     0 |       97696 |             5 |      14352 |          12 |      83344
 TopTransactionContext |             | TopMemoryContext     |     1 |        8192 |             1 |       7760 |           0 |        432
 Btree proof lookup cache |          | TopMemoryContext     |     1 |        8192 |             1 |        576 |           0 |       7616
 TableSpace cache      |             | TopMemoryContext     |     1 |        8192 |             1 |       2112 |           0 |       6080
 Type information cache |            | TopMemoryContext     |     1 |       24384 |             2 |       2640 |           0 |      21744
 Operator lookup cache |             | TopMemoryContext     |     1 |       24576 |             2 |      10776 |           3 |      13800
 Record information cache |          | TopMemoryContext     |     1 |        8192 |             1 |       1600 |           0 |       6592
 RowDescriptionContext |             | TopMemoryContext     |     1 |        8192 |             1 |       6912 |           0 |       1280
 MessageContext        |             | TopMemoryContext     |     1 |       65536 |             4 |      32720 |           2 |      32816
 search_path processing cache |      | TopMemoryContext     |     1 |        8192 |             1 |       5616 |           8 |       2576
 Operator class cache  |             | TopMemoryContext     |     1 |        8192 |             1 |        576 |           0 |       7616
 PgStat Shared Ref Hash |            | TopMemoryContext     |     2 |        7232 |             2 |        704 |           0 |       6528
 PgStat Shared Ref     |             | TopMemoryContext     |     1 |        8192 |             4 |       4072 |           2 |       4120
 PgStat Pending        |             | TopMemoryContext     |     1 |       16384 |             5 |      15984 |          48 |        400
 smgr relation table   |             | TopMemoryContext     |     1 |       32768 |             3 |      16848 |           8 |      15920
 TransactionAbortContext |           | TopMemoryContext     |     1 |       32768 |             1 |      32528 |           0 |        240
 Portal hash           |             | TopMemoryContext     |     1 |        8192 |             1 |        576 |           0 |       7616
 TopPortalContext      |             | TopMemoryContext     |     1 |        8192 |             1 |       7680 |           0 |        512
 PortalContext         | <unnamed>   | TopPortalContext     |     2 |        1024 |             1 |        592 |           0 |        432
 ExecutorState         |             | PortalContext        |     3 |       49216 |             4 |      13424 |           1 |      35792
(20 rows)
```

# Extensions & Tools

- pg_buffercache

- pg_prewarm

- pmap (Linux)

```
00005583f3a14000    220K rw---   [ anon ]
00005583f459c000    924K rw---   [ anon ]
00005583f4683000   1392K rw---   [ anon ]
00007f877932f000    776K rw---   [ anon ]
00007f87793f1000   1024K rw-s- /dev/shm/PostgreSQL.4106567874
00007f87794f1000    132K rw---   [ anon ]
00007f8779512000     28K rw-s- /dev/shm/PostgreSQL.1502130758
00007f8779519000 146336K rw-s- /dev/zero (deleted)
00007f8782401000    348K r---- /usr/lib/locale/C.utf8/LC_CTYPE
00007f8782458000     28K r--s- /usr/lib/x86_64-linux-gnu/gconv/gconv-modules.cache
00007f878245f000     20K rw---   [ anon ]
00007f8782464000    160K r---- /usr/lib/x86_64-linux-gnu/libc.so.6
00007f878248c000   1620K r-x-- /usr/lib/x86_64-linux-gnu/libc.so.6
00007f8782621000    352K r---- /usr/lib/x86_64-linux-gnu/libc.so.6
00007f8782679000      4K ----- /usr/lib/x86_64-linux-gnu/libc.so.6
00007f878267a000     16K r---- /usr/lib/x86_64-linux-gnu/libc.so.6
00007f878267e000      8K rw--- /usr/lib/x86_64-linux-gnu/libc.so.6
00007f8782680000     52K rw---   [ anon ]
00007f878268d000      8K r---- /usr/lib/x86_64-linux-gnu/libz.so.1.2.11
00007f878268f000     68K r-x-- /usr/lib/x86_64-linux-gnu/libz.so.1.2.11
00007f878269a000     24K r---- /usr/lib/x86_64-linux-gnu/libz.so.1.2.11
00007f87826a6000      4K ----- /usr/lib/x86_64-linux-gnu/libz.so.1.2.11
00007f87826a7000      4K r---- /usr/lib/x86_64-linux-gnu/libz.so.1.2.11
00007f87826a8000      4K rw--- /usr/lib/x86_64-linux-gnu/libz.so.1.2.11
00007f87826a9000    712K r---- /usr/lib/x86_64-linux-gnu/libcrypto.so.3
00007f878275b000   2424K r-x-- /usr/lib/x86_64-linux-gnu/libcrypto.so.3
00007f87829b9000    840K r---- /usr/lib/x86_64-linux-gnu/libcrypto.so.3
00007f8782a8b000    364K r---- /usr/lib/x86_64-linux-gnu/libcrypto.so.3
00007f8782ae6000     12K rw--- /usr/lib/x86_64-linux-gnu/libcrypto.so.3
00007f8782ae9000     12K rw---   [ anon ]
00007f8782aec000    120K r---- /usr/lib/x86_64-linux-gnu/libssl.so.3
00007f8782b0a000    364K r-x-- /usr/lib/x86_64-linux-gnu/libssl.so.3
00007f8782b65000    116K r---- /usr/lib/x86_64-linux-gnu/libssl.so.3
00007f8782b82000     40K r---- /usr/lib/x86_64-linux-gnu/libssl.so.3
00007f8782b8c000     16K rw--- /usr/lib/x86_64-linux-gnu/libssl.so.3
00007f8782b90000     56K r---- /usr/lib/x86_64-linux-gnu/libm.so.6
00007f8782b9e000    496K r-x-- /usr/lib/x86_64-linux-gnu/libm.so.6
00007f8782c1a000    364K r---- /usr/lib/x86_64-linux-gnu/libm.so.6
00007f8782c75000      4K r---- /usr/lib/x86_64-linux-gnu/libm.so.6
00007f8782c76000      4K rw--- /usr/lib/x86_64-linux-gnu/libm.so.6
00007f8782c78000      4K rw-s-   [ shmid=0x0 ]
00007f8782c79000      4K r---- /usr/lib/locale/C.utf8/LC_TIME
00007f8782c7a000      4K r---- /usr/lib/locale/C.utf8/LC_NUMERIC
00007f8782c7b000      4K r---- /usr/lib/locale/C.utf8/LC_MONETARY
00007f8782c7c000      4K r---- /usr/lib/locale/C.utf8/LC_MESSAGES/SYS_LC_MESSAGES
00007f8782c7d000      8K rw---   [ anon ]
00007f8782c7f000      8K r---- /usr/lib/x86_64-linux-gnu/ld-linux-x86-64.so.2
00007f8782c81000    168K r-x-- /usr/lib/x86_64-linux-gnu/ld-linux-x86-64.so.2
00007f8782cab000     44K r---- /usr/lib/x86_64-linux-gnu/ld-linux-x86-64.so.2
00007f8782cb6000      4K r---- /usr/lib/locale/C.utf8/LC_COLLATE
00007f8782cb7000      8K r---- /usr/lib/x86_64-linux-gnu/ld-linux-x86-64.so.2
00007f8782cb9000      8K rw--- /usr/lib/x86_64-linux-gnu/ld-linux-x86-64.so.2
00007ffe0ec66000    136K rw---   [ stack ]
00007ffe0ed2e000     16K r----   [ anon ]
00007ffe0ed32000      8K r-x--   [ anon ]
```

# Projects

## Ideas/unmerged

- Invalidate buffer cache – patch is out
- Memory shrink/expand – serverless, overbooking (anyone ? ☺)
- Memory accounting & limiting
- Merge SLRUs into buffer pool
- RelCache cleanups.
- Improve allocation speeds
- New buffer pool replacement algorithm.

# Acknowledgements

- Microsoft Contributors/Committers team.
- Specially thanks - Thomas, David, Andres & Teresa.