# DBMS Implementation : InnoDB Initialize

woonhak.kang

woonagi319@skku.edu

VLDB Lab. @ SKKU
http://vldb.skku.ac.kr

# Contents

- Before we begin : tools
  - message logging and entry point
  - Coding conventions : the very first

- InnoDB Initialize
  - Init. in-memory objects and check configurations
  - aio system, file, buffer pool, log system

- Create objects
  - lock-system, io-handler, data files and log files

- Open
  - log and system table space
  - transaction system

- Recovery if necessary

- Create double write system if necessary

# Message Log for development

- ib_logf(log_level, "message")
  - or you can also do fprintf(stderr, "")

- log levels
  - in file : include/ha_prototypes.h:449

```
enum ib_log_level_t {
    IB_LOG_LEVEL_INFO,
    IB_LOG_LEVEL_WARN,
    IB_LOG_LEVEL_ERROR,
    IB_LOG_LEVEL_FATAL
};
```

# InnoDB entry

- MYSQL_SRC
  - e.g.) /home/woonhak/workspace/mysql-5.6.21

- INNODB_SRC
  - $MYSQL_SRC/storage/innobase/

- InnoDB entry point
  - in file: srv/srv0start.cc:1502 innobase_start_or_create_for_mysql(void)
  - Starts InnoDB and creates a new database if database files are not found and the user wants.

# Coding conventions

- srv_XXX : global system variables

- end of #if~#else~#endif directive

- module/function name

```
1449 #ifdef UNIV_SYNC_DEBUG
1450   /* Create the thread latch level array where the latch levels
1451   are stored for each OS thread */
1452
1453   sync_thread_level_arrays = static_cast<sync_thread_t*>(
1454     calloc(sizeof(sync_thread_t), OS_THREAD_MAX_N));
1455
1456   ut_a(sync_thread_level_arrays != NULL);
1457
1458 #endif /* UNIV_SYNC_DEBUG */
```

  - use abbr.
    - for e.g.)
    - dictionary -> dict, server->srv, utility->ut, memory->mem, buffer->buf

- Assertion
  - ut_a(), ut_ad()

# InnoDB Initialize

- We are under
  - Linux
  - native aio (libaio)
  - use O_DIRECT
  - does not use raw_disk (files on file system)
  - default database files are already created

- Initialize
  - aio system
  - file system
  - buffer pool
  - log system

# InnoDB Initialize

- srv/srv0start.cc:1502

```
1496 /***********************************************************
1497 Starts InnoDB and creates a new database if database files
1498 are not found and the user wants.
1499 @return DB_SUCCESS or error code */
1500 UNIV_INTERN
1501 dberr_t
1502 innobase_start_or_create_for_mysql(void)
1503 /*====================================*/
1504 {
1505   ibool    create_new_db;
1506   lsn_t    min_flushed_lsn;
1507   lsn_t    max_flushed_lsn;
1508 #ifdef UNIV_LOG_ARCHIVE
1509   ulint    min_arch_log_no;
1510   ulint    max_arch_log_no;
1511 #endif /* UNIV_LOG_ARCHIVE */
1512   ulint    sum_of_new_sizes;
1513   ulint    sum_of_data_file_sizes;
1514   ulint    tablespace_size_in_header;
1515   dberr_t    err;
1516   unsigned  i;
1517   ulint    srv_n_log_files_found = srv_n_log_files;
1518   ulint    io_limit;
1519   mtr_t    mtr;
1520   ib_bh_t*  ib_bh;
1521   ulint    n_recovered_trx;
1522   char    logfilename[10000];
1523   char*    logfile0  = NULL;
1524   size_t     dirnamelen;
1525
1526   if (srv_force_recovery > SRV_FORCE_NO_TRX_UNDO) {
1527     srv_read_only_mode = true;
1528   }
1529
1530   if (srv_read_only_mode) {
1531     ib_logf(IB_LOG_LEVEL_INFO, "Started in read only mode");
1532   }
```

# InnoDB Initialize

- srv_use_native_aio
  - include/srv0srv.cc

```
142 /* If this flag is TRUE, then we will use the native aio of the
143 OS (provided we compiled Innobase with it in), otherwise we will
144 use simulated aio we build below with threads.
145 Currently we support native aio on windows and linux */
146 UNIV_INTERN my_bool srv_use_native_aio = TRUE;
```

- srv/srv0start.cc:17xx
  - Check use native AIO

```
1728 #elif defined(LINUX_NATIVE_AIO)
1729
1730   if (srv_use_native_aio) {
1731     ib_logf(IB_LOG_LEVEL_INFO, "Using Linux native AIO");
1732   }
1733 #else
1734   /* Currently native AIO is supported only on windows and linux
1735   and that also when the support is compiled in. In all other
1736   cases, we ignore the setting of innodb_use_native_aio. */
1737   srv_use_native_aio = FALSE;
1738 #endif /* __WIN__ */
```

# InnoDB Initialize

- `srv/srv0start.cc:17XX`

- file flush method : `srv_unix_file_flush_method`
  - `fsync` : use fsync() to flush both log and data - default
  - `o_dsync` : use O_SYNC for log, fsync() for data
  - `o_direct` : use O_DIRECT for data, log is buffered, use fsync() to flush log and data
  - `o_direct_no_fsync` : use O_DIRECT but skip fsync for data

```
1740  if (srv_file_flush_method_str == NULL) {
1741    /* These are the default options */
1742
1743    srv_unix_file_flush_method = SRV_UNIX_FSYNC;
1744
1745    srv_win_file_flush_method = SRV_WIN_IO_UNBUFFERED;
1746 #ifndef __WIN__
1747  } else if (0 == ut_strcmp(srv_file_flush_method_str, "fsync")) {
1748    srv_unix_file_flush_method = SRV_UNIX_FSYNC;
1749
1750  } else if (0 == ut_strcmp(srv_file_flush_method_str, "O_DSYNC")) {
1751    srv_unix_file_flush_method = SRV_UNIX_O_DSYNC;
1752
1753  } else if (0 == ut_strcmp(srv_file_flush_method_str, "O_DIRECT")) {
1754    srv_unix_file_flush_method = SRV_UNIX_O_DIRECT;
1755
1756  } else if (0 == ut_strcmp(srv_file_flush_method_str, "O_DIRECT_NO_FSYNC")) {
1757    srv_unix_file_flush_method = SRV_UNIX_O_DIRECT_NO_FSYNC;
```

# buffer pool size threshold

- `srv/srv0start.cc`
  - if buffer pool size less than threshold than buf pool instances = 1
  - `#define BUF_POOL_SIZE_THRESHOLD (1024 * 1024 * 1024)`

```
1792  #define BUF_POOL_SIZE_THRESHOLD (1024 * 1024 * 1024)
1793    srv_max_n_threads = 1    /* io_ibuf_thread */
1794            + 1 /* io_log_thread */
1795            + 1 /* lock_wait_timeout_thread */
1796            + 1 /* srv_error_monitor_thread */
1797            + 1 /* srv_monitor_thread */
1798            + 1 /* srv_master_thread */
1799            + 1 /* srv_purge_coordinator_thread */
1800            + 1 /* buf_dump_thread */
1801            + 1 /* dict_stats_thread */
1802            + 1 /* fts_optimize_thread */
1803            + 1 /* recv_writer_thread */
1804            + 1 /* buf_flush_page_cleaner_thread */
1805            + 1 /* trx_rollback_or_clean_all_recovered */
1806            + 128 /* added as margin, for use of
1807            InnoDB Memcached etc. */
1808            + max_connections
1809            + srv_n_read_io_threads
1810            + srv_n_write_io_threads
1811            + srv_n_purge_threads
1812            /* FTS Parallel Sort */
1813            + fts_sort_pll_degree * FTS_NUM_AUX_INDEX
1814              * max_connections;
1815
1816    if (srv_buf_pool_size < BUF_POOL_SIZE_THRESHOLD) {
1817      /* If buffer pool is less than 1 GB,
1818      use only one buffer pool instance */
1819      srv_buf_pool_instances = 1;
1820    }
```

# server boot

- Boots the InnoDB server
  - invokes `srv_boot()` at `srv0start.cc:1822`

- `srv/srv0srv.cc:1057,srv_boot()`

```
1053 /***************************************************************//**
1054 Boots the InnoDB server. */
1055 UNIV_INTERN
1056 void
1057 srv_boot(void)
1058 /*==========*/
1059 {
1060     /* Transform the init parameter values given by MySQL to
1061     use units we use inside InnoDB: */
1062
1063     srv_normalize_init_values();
1064
1065     /* Initialize synchronization primitives, memory management, and thread
1066     local storage */
1067
1068     srv_general_init();
1069
1070     /* Initialize this module */
1071
1072     srv_init();
1073     srv_mon_create();
1074 }
```

init general variables and objects

# server boot

- `srv/srv0srv.cc:1013, srv_general_init()`

```
1008 /**********************************************************//**
1009 Initializes the synchronization primitives, memory system, and the thread
1010 local storage. */
1011 UNIV_INTERN
1012 void
1013 srv_general_init(void)
1014 /*==================*/
1015 {
1016   ut_mem_init();
1017   /* Reset the system variables in the recovery module. */
1018   recv_sys_var_init();
1019   os_sync_init();
1020   sync_init();
1021   mem_init(srv_mem_pool_size);
1022   que_init();
1023   row_mysql_init();
1024 }
```

os_sync_mutex init

mutex init

additional memory pool (not buffer pool)

# server boot

- ## srv/srv0srv.cc:1057, srv_boot()

```
1053 /*******************************************************************//**
1054 Boots the InnoDB server. */
1055 UNIV_INTERN
1056 void
1057 srv_boot(void)
1058 /*==========*/
1059 {
1060   /* Transform the init parameter values given by MySQL to
1061   use units we use inside InnoDB: */
1062
1063   srv_normalize_init_values();
1064
1065   /* Initialize synchronization primitives, memory management, and thread
1066   local storage */
1067
1068   srv_general_init();
1069
1070   /* Initialize this module */
1071
1072   srv_init();
1073   srv_mon_create();
1074 }
```

server init.

# server boot

- `srv/srv0srv.cc:909, srv_init()`

```
905  /******************************************************************//**
906  Initializes the server. */
907  UNIV_INTERN
908  void
909  srv_init(void)
910  /*==========*/
911  {
912    ulint n_sys_threads = 0;
913    ulint srv_sys_sz = sizeof(*srv_sys);
914
915  #ifndef HAVE_ATOMIC_BUILTINS
916    mutex_create(server_mutex_key, &server_mutex, SYNC_ANY_LATCH);
917  #endif /* !HAVE_ATOMIC_BUILTINS */
918
919    mutex_create(srv_innodb_monitor_mutex_key,
920          &srv_innodb_monitor_mutex, SYNC_NO_ORDER_CHECK);
921
922    if (!srv_read_only_mode) {
923
924      /* Number of purge threads + master thread */
925      n_sys_threads = srv_n_purge_threads + 1;
926
927      srv_sys_sz += n_sys_threads * sizeof(*srv_sys->sys_threads);
928    }
929
930    srv_sys = static_cast<srv_sys_t*>(mem_zalloc(srv_sys_sz));
931
932    srv_sys->n_sys_threads = n_sys_threads;
```

init. system and system threads

# of system threads = purge thread + master thread

# server boot

- `srv/srv0srv.cc:909, srv_init()`

```
934   if (!srv_read_only_mode) {
935
936     mutex_create(srv_sys_mutex_key, &srv_sys->mutex, SYNC_THREADS);
937
938     mutex_create(srv_sys_tasks_mutex_key,
939             &srv_sys->tasks_mutex, SYNC_ANY_LATCH);
940
941     srv_sys->sys_threads = (srv_slot_t*) &srv_sys[1];
942
943     for (ulint i = 0; i < srv_sys->n_sys_threads; ++i) {
944       srv_slot_t* slot = &srv_sys->sys_threads[i];
945
946       slot->event = os_event_create();
947
948       ut_a(slot->event);
949     }
950
951     srv_error_event = os_event_create();
952
953     srv_monitor_event = os_event_create();
954
955     srv_buf_dump_event = os_event_create();
956
957     UT_LIST_INIT(srv_sys->tasks);
958   }
```

mutex for srv_sys

system threads

create event for synchronization and signaling

# server boot

- `srv/srv0srv.cc:909, srv_init()`

```
972    /* Create dummy indexes for infimum and supremu
973
974    dict_ind_init();           ←  dictionary index
975
976    srv_conc_init();
977
978    /* Initialize some INFORMATION SCHEMA internal structures */
979    trx_i_s_cache_init(trx_i_s_cache);
980
981    ut_crc32_init();
982
983    dict_mem_init();
984 }
```

# server boot

- `dict/dict0dict.cc:5759, dict_ind_init()`

```
5755 /*******************************************************************//**
5756 Inits dict_ind_redundant and dict_ind_compact.
5757 UNIV_INTERN
5758 void
5759 dict_ind_init(void)
5760 /*===============*/
5761 {
5762   dict_table_t*   table;
5763
5764   /* create dummy table and index for REDUNDANT infimum and supremum */
5765   table = dict_mem_table_create("SYS_DUMMY1", DICT_HDR_SPACE, 1, 0, 0);
5766   dict_mem_table_add_col(table, NULL, NULL, DATA_CHAR,
5767             DATA_ENGLISH | DATA_NOT_NULL, 8);
5768
5769   dict_ind_redundant = dict_mem_index_create("SYS_DUMMY1", "SYS_DUMMY1",
5770             DICT_HDR_SPACE, 0, 1);
5771   dict_index_add_col(dict_ind_redundant, table,
5772         dict_table_get_nth_col(table, 0), 0);
5773   dict_ind_redundant->table = table;
5774
5775   /* create dummy table and index for COMPACT infimum and supremum */
5776   table = dict_mem_table_create("SYS_DUMMY2",
5777             DICT_HDR_SPACE, 1,
5778             DICT_TF_COMPACT, 0);
5779   dict_mem_table_add_col(table, NULL, NULL, DATA_CHAR,
5780             DATA_ENGLISH | DATA_NOT_NULL, 8);
5781   dict_ind_compact = dict_mem_index_create("SYS_DUMMY2", "SYS_DUMMY2",
5782             DICT_HDR_SPACE, 0, 1);
5783   dict_index_add_col(dict_ind_compact, table,
5784         dict_table_get_nth_col(table, 0), 0);
5785   dict_ind_compact->table = table;
5786
5787   /* avoid ut_ad(index->cached) in dict_index_get_n_unique_in_tree */
5788   dict_ind_redundant->cached = dict_ind_compact->cached = TRUE;
5789 }
```

dictionary table

by default : we use
compact type

# server boot

- `srv/srv0srv.cc:909, srv_init()`

```
972    /* Create dummy indexes for infimum and supremum records */
973
974    dict_ind_init();
975
976    srv_conc_init();
977
978    /* Initialize some INFORMATION SCHEMA internal structures */
979    trx_i_s_cache_init(trx_i_s_cache);
980
981    ut_crc32_init();
982
983    dict_mem_init();
984 }
```

checksum utility

# server boot

- ut/ut0crc32.cc:276, ut_crc32_init()

```
271 /*************************************************************************//**
272 Initializes the data structures used by ut_crc32(). Does not do any
273 allocations, would not hurt if called twice, but would be pointless. */
274 UNIV_INTERN
275 void
276 ut_crc32_init()
277 /*===========*/
278 {
279 #if defined(__GNUC__) && defined(__x86_64__)
280   ib_uint32_t vend[3];
281   ib_uint32_t model;
282   ib_uint32_t family;
283   ib_uint32_t stepping;
284   ib_uint32_t features_ecx;
285   ib_uint32_t features_edx;
286
287   ut_cpuid(vend, &model, &family, &stepping,
288       &features_ecx, &features_edx);
312   if (ut_crc32_sse2_enabled) {
313     ut_crc32 = ut_crc32_sse42;
314   } else {
315     ut_crc32_slice8_table_init();
316     ut_crc32 = ut_crc32_slice8;
317   }
318 }
```

for checksum - use crc32 sse inst.

# InnoDB Initialize

- `srv/srv0start.cc`

- total io threads

```
1883    /* If user has set the value of innodb_file_io_threads then
1884    we'll emit a message telling the user that this parameter
1885    is now deprecated. */
1886    if (srv_n_file_io_threads != 4) {
1887      ib_logf(IB_LOG_LEVEL_WARN,
1888        "innodb_file_io_threads is deprecated. Please use "
1889        "innodb_read_io_threads and innodb_write_io_threads "
1890        "instead");
1891    }
1892
1893    /* Now overwrite the value on srv_n_file_io_threads */
1894    srv_n_file_io_threads = srv_n_read_io_threads;
1895
1896    if (!srv_read_only_mode) {
1897      /* Add the log and ibuf IO threads. */
1898      srv_n_file_io_threads += 2;
1899      srv_n_file_io_threads += srv_n_write_io_threads;
1900    } else {
1901      ib_logf(IB_LOG_LEVEL_INFO,
1902        "Disabling background IO write threads.");
1903
1904      srv_n_write_io_threads = 0;
1905    }
```

# of io threads = read + write + 2

# InnoDB Initialize

- `srv/srv0start.cc`

- aio system init : os_aio_init()

```
1907    ut_a(srv_n_file_io_threads <= SRV_MAX_N_IO_THREADS);
1908
1909    io_limit = 8 * SRV_N_PENDING_IOS_PER_THREAD;
1910
1911    /* On Windows when using native aio the number of aio requests
1912    that a thread can handle at a given time is limited to 32
1913    i.e.: SRV_N_PENDING_IOS_PER_THREAD */
1914 # ifdef __WIN__
1915    if (srv_use_native_aio) {
1916      io_limit = SRV_N_PENDING_IOS_PER_THREAD;
1917    }
1918 # endif /* __WIN__ */
1919
1920    if (!os_aio_init(io_limit,
1921        srv_n_read_io_threads,
1922        srv_n_write_io_threads,
1923        SRV_MAX_N_PENDING_SYNC_IOS)) {
1924
1925      ib_logf(IB_LOG_LEVEL_ERROR,
1926        "Fatal : Cannot initialize AIO sub-system");
1927
1928      return(DB_ERROR);
1929    }
```
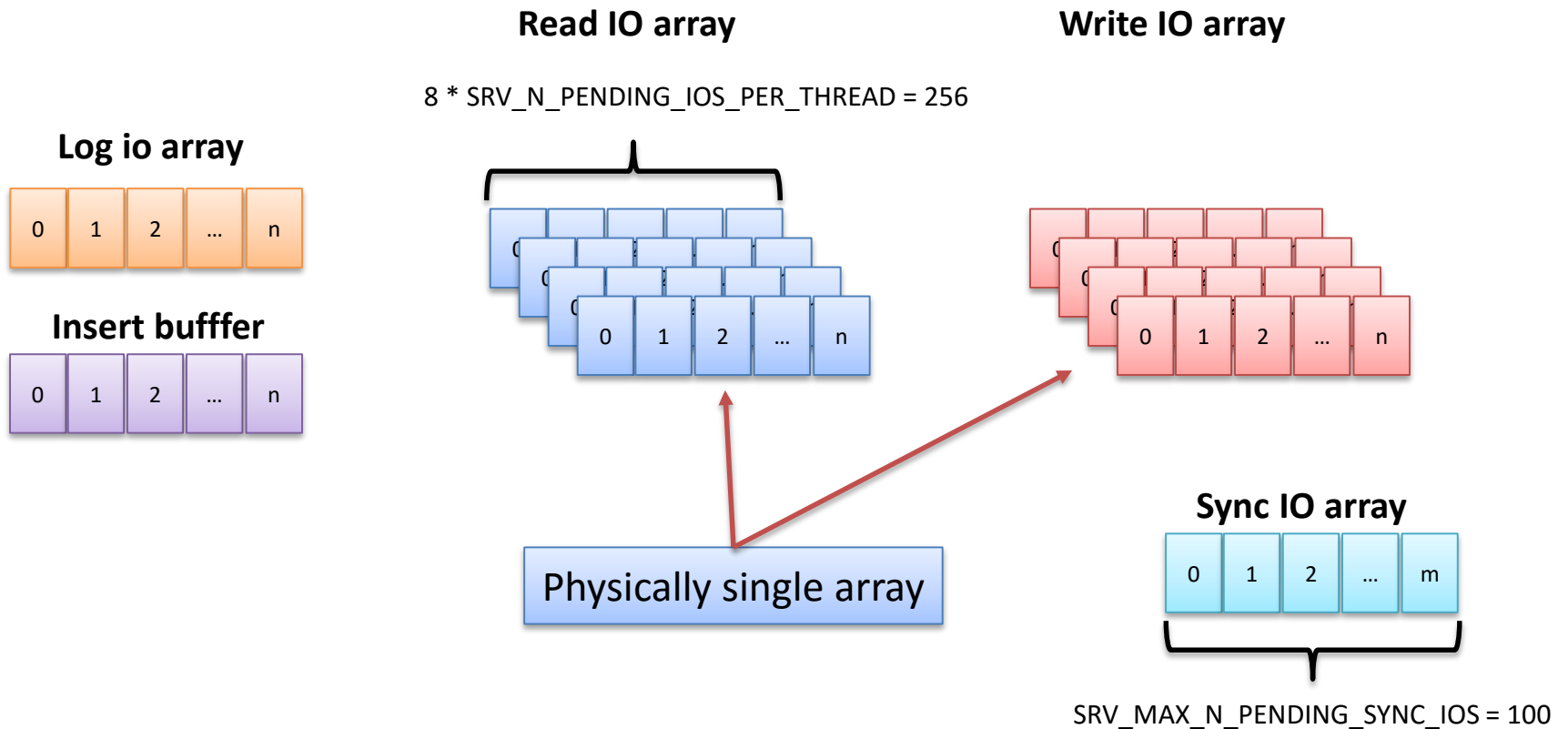
n_pending_ios_per_thread = 32

aio system init

# aio system init

- Initializes the asynchronous io system.

- Creates one array each for ibuf and log i/o.

- Also creates one array each for read and write where each array is divided logically into n_read_segs and n_write_segs respectively.

- The caller must create an i/o handler thread for each segment in these arrays.

- This function also creates the sync array.

- No i/o handler thread needs to be created for that

# aio system init

- aio thread and array



Read IO array

8 * SRV_N_PENDING_IOS_PER_THREAD = 256

Write IO array

Log io array

Insert bufffer

Physically single array

Sync IO array

SRV_MAX_N_PENDING_SYNC_IOS = 100

# aio system init

- `os/os0file.cc:3835, os_aio_init()`

```
3833 UNIV_INTERN
3834 ibool
3835 os_aio_init(
3836 /*========*/
3837   ulint n_per_seg,   /*<! in: maximum number of pending aio
3838         operations allowed per segment */
3839   ulint n_read_segs,  /*<! in: number of reader threads */
3840   ulint n_write_segs, /*<! in: number of writer threads */
3841   ulint n_slots_sync) /*<! in: number of slots in the sync aio
3842         array */
3843 {
3844   os_io_init_simple();
3845
3846 #if defined(LINUX_NATIVE_AIO)
3847   /* Check if native aio is supported on this system and tmpfs */
3848   if (srv_use_native_aio && !os_aio_native_aio_supported()) {
3849
3850     ib_logf(IB_LOG_LEVEL_WARN, "Linux Native AIO disabled.");
3851
3852     srv_use_native_aio = FALSE;
3853   }
3854 #endif /* LINUX_NATIVE_AIO */
```

check aio *really* supported

# aio system init

- `os/os0file.cc:3578,`
  os_aio_native_aio_supported()
  – make a temp io context
  – and do test read or write

```
3578 os_aio_native_aio_supported(void)
3579 /*==============================*/
3580 {
3581   int       fd;
3582   io_context_t     io_ctx;
3583   char       name[1000];
3584
3585   if (!os_aio_linux_create_io_ctx(1, &io_ctx)) {
3586     /* The platform does not support native aio. */
3587     return(FALSE);
3588   } else if (!srv_read_only_mode) {
3589     /* Now check if tmpdir supports native aio ops. */
3590     fd = innobase_mysql_tmpfile();
3591
3592     if (fd < 0) {
3593       ib_logf(IB_LOG_LEVEL_WARN,
3594         "Unable to create temp file to check "
3595         "native AIO support.");
3596
3597       return(FALSE);
3598     }
3599   } else {
```

# aio system init

- `os/os0file.cc:3835, os_aio_init()`

- prepare read io threads

make aio array (kind of io queue) see it later

```
3858    os_aio_read_array = os_aio_array_create(
3859      n_read_segs * n_per_seg, n_read_segs);
3860
3861    if (os_aio_read_array == NULL) {
3862      return(FALSE);
3863    }
3864
3865    ulint start = (srv_read_only_mode) ? 0 : 2;
3866    ulint n_segs = n_read_segs + start;
3867
3868    /* 0 is the ibuf segment and 1 is the insert buffer segment. */
3869    for (ulint i = start; i < n_segs; ++i) {
3870      ut_a(i < SRV_MAX_N_IO_THREADS);
3871      srv_io_thread_function[i] = "read thread";
3872    }
```

# aio system init

- `os/os0file.cc:3835, os_aio_init()`

- prepare log and ibuf thread

```
3876    if (!srv_read_only_mode) {
3877
3878      os_aio_log_array = os_aio_array_create(n_per_seg, 1);
3879
3880      if (os_aio_log_array == NULL) {
3881        return(FALSE);
3882      }
3883
3884      ++n_segments;
3885
3886      srv_io_thread_function[1] = "log thread";
3887
3888      os_aio_ibuf_array = os_aio_array_create(n_per_seg, 1);
3889
3890      if (os_aio_ibuf_array == NULL) {
3891        return(FALSE);
3892      }
3893
3894      ++n_segments;
3895
3896      srv_io_thread_function[0] = "insert buffer thread";
3897
3898      os_aio_write_array = os_aio_array_create(
3899        n_write_segs * n_per_seg, n_write_segs);
3900
3901      if (os_aio_write_array == NULL) {
3902        return(FALSE);
3903      }
```

# aio system init

- `os/os0file.cc:3835, os_aio_init()`

- prepare write io thread and sync array

```
3898     os_aio_write_array = os_aio_array_create(
3899       n_write_segs * n_per_seg, n_write_segs);
3900
3901     if (os_aio_write_array == NULL) {
3902       return(FALSE);
3903     }
3904
3905     n_segments += n_write_segs;
3906
3907     for (ulint i = start + n_read_segs; i < n_segments; ++i) {
3908       ut_a(i < SRV_MAX_N_IO_THREADS);
3909       srv_io_thread_function[i] = "write thread";
3910     }
3911
3912     ut_ad(n_segments >= 4);
3913   } else {
3914     ut_ad(n_segments > 0);
3915   }
3916
3917   os_aio_sync_array = os_aio_array_create(n_slots_sync, 1);
3918
3919   if (os_aio_sync_array == NULL) {
3920     return(FALSE);
3921   }
```

# aio system init

- `os/os0file.cc:3689, create aio array`

```
3689 os_aio_array_create(
3690 /*=================*/
3691   ulint n,      /*!< in: maximum number of pending aio
3692          operations allowed; n must be
3693          divisible by n_segments */
3694   ulint n_segments) /*!< in: number of segments in the aio array */
3695 {
3696   os_aio_array_t* array;
3697 #ifdef WIN_ASYNC_IO
3698   OVERLAPPED* over;
3699 #elif defined(LINUX_NATIVE_AIO)
3700   struct io_event*  io_event = NULL;
3701 #endif /* WIN_ASYNC_IO */
3702   ut_a(n > 0);
3703   ut_a(n_segments > 0);
3704
3705   array = static_cast<os_aio_array_t*>(ut_malloc(sizeof(*array)));
3706   memset(array, 0x0, sizeof(*array));
3707
3708   array->mutex = os_mutex_create();
3709   array->not_full = os_event_create();
3710   array->is_empty = os_event_create();
3711
3712   os_event_set(array->is_empty);
3713
3714   array->n_slots = n;
3715   array->n_segments = n_segments;
3716
3717   array->slots = static_cast<os_aio_slot_t*>(
3718     ut_malloc(n * sizeof(*array->slots)));
3719
3720   memset(array->slots, 0x0, sizeof(n * sizeof(*array->slots)));
```

create array struct

create mutex for the array

create slots

# aio system init

- `os/os0file.cc:3689, create aio array`

```
3725 #if defined(LINUX_NATIVE_AIO)
3726   array->aio_ctx = NULL;
3727   array->aio_events = NULL;
3728
3729   /* If we are not using native aio interface then skip this
3730   part of initialization. */
3731   if (!srv_use_native_aio) {
3732     goto skip_native_aio;
3733   }
3734
3735   /* Initialize the io_context array. One io_context
3736   per segment in the array. */
3737
3738   array->aio_ctx = static_cast<io_context**>(
3739     ut_malloc(n_segments * sizeof(*array->aio_ctx)));
3740
3741   for (ulint i = 0; i < n_segments; ++i) {
3742     if (!os_aio_linux_create_io_ctx(n/n_segments,
3743            &array->aio_ctx[i])) {
3744       /* If something bad happened during aio setup
3745       we should call it a day and return right away.
3746       We don't care about any leaks because a failure
3747       to initialize the io subsystem means that the
3748       server (or atleast the innodb storage engine)
3749       is not going to startup. */
3750       return(NULL);
3751     }
3752   }
3753
3754   /* Initialize the event array. One event per slot. */
3755   io_event = static_cast<struct io_event*>(
3756     ut_malloc(n * sizeof(*io_event)));
3757
3758   memset(io_event, 0x0, sizeof(*io_event) * n);
3759   array->aio_events = io_event;
```

io context

io event array

# aio system init

- os_aio_array_t

```
190 /** The asynchronous i/o array structure */
191 struct os_aio_array_t{
192   os_ib_mutex_t mutex;   /*!< the mutex protecting the aio array */
193   os_event_t  not_full;
194        /*!< The event which is set to the
195        signaled state when there is space in
196        the aio outside the ibuf segment */
197   os_event_t  is_empty;
198        /*!< The event which is set to the
199        signaled state when there are no
200        pending i/os in this array */
201   ulint   n_slots;/*!< Total number of slots in the aio
202        array.  This must be divisible by
203        n_threads. */
204   ulint   n_segments;
205        /*!< Number of segments in the aio
206        array of pending aio requests. A
207        thread can wait separately for any one
208        of the segments. */
209   ulint   cur_seg;/*!< We reserve IO requests in round
210        robin fashion to different segments.
211        This points to the segment that is to
212        be used to service next IO request. */
213   ulint   n_reserved;
214        /*!< Number of reserved slots in the
215        aio array outside the ibuf segment */
216   os_aio_slot_t*  slots;  /*!< Pointer to the slots in the array */
```

```
227 #if defined(LINUX_NATIVE_AIO)
228   io_context_t*    aio_ctx;
229        /* completion queue for IO. There is
230        one such queue per segment. Each thread
231        will work on one ctx exclusively. */
232   struct io_event*  aio_events;
233        /* The array to collect completed IOs.
234        There is one such event for each
235        possible pending IO. The size of the
236        array is equal to n_slots. */
237 #endif /* LINUX_NATIV_AIO */
238 };
```

# aio system init

- os_aio_slot_t

```
154 /** The asynchronous i/o array slot structure */
155 struct os_aio_slot_t{
156   ibool   is_read;  /*!< TRUE if a read operation */
157   ulint   pos;      /*!< index of the slot in the aio
158           array */
159   ibool   reserved; /*!< TRUE if this slot is reserved */
160   time_t    reservation_time;/*!< time when reserved */
161   ulint   len;      /*!< length of the block to read or
162           write */
163   byte*   buf;      /*!< buffer used in i/o */
164   ulint   type;     /*!< OS_FILE_READ or OS_FILE_WRITE */
165   os_offset_t offset;   /*!< file offset in bytes */
166   os_file_t file;   /*!< file where to read or write */
167   const char* name;   /*!< file name or path */
168   ibool   io_already_done;/*!< used only in simulated aio:
169           TRUE if the physical i/o already
170           made and only the slot message
171           needs to be passed to the caller
172           of os_aio_simulated_handle */
173   fil_node_t* message1; /*!< message which is given by the */
174   void*   message2; /*!< the requester of an aio operation
175           and which can be used to identify
176           which pending aio operation was
177           completed */
183 #elif defined(LINUX_NATIVE_AIO)
184   struct iocb control;  /* Linux control block for aio */
185   int   n_bytes;  /* bytes written/read. */
186   int   ret;    /* AIO return code */
187 #endif /* WIN_ASYNC_IO */
188 };
```

# InnoDB Initialize

- file system init

- `srv0start.cc:1931, fil_init()`

invoke file system init

```
1931    fil_init(srv_file_per_table ? 50000 : 5000, srv_max_n_open_files);
1932
1933    double  size;
1934    char    unit;
1935
1936    if (srv_buf_pool_size >= 1024 * 1024 * 1024) {
1937      size = ((double) srv_buf_pool_size) / (1024 * 1024 * 1024);
1938      unit = 'G';
1939    } else {
1940      size = ((double) srv_buf_pool_size) / (1024 * 1024);
1941      unit = 'M';
1942    }
1943
1944    /* Print time to initialize the buffer pool */
1945    ib_logf(IB_LOG_LEVEL_INFO,
1946      "Initializing buffer pool, size = %.1f%c", size, unit);
1947
1948    err = buf_pool_init(srv_buf_pool_size, srv_buf_pool_instances);
```

# file system init

- ## fil_system_t

```
261 struct fil_system_t {
262 #ifndef UNIV_HOTBACKUP
263   ib_mutex_t    mutex;    /*!< The mutex protecting the cache */
264 #endif /* !UNIV_HOTBACKUP */
265   hash_table_t* spaces;   /*!< The hash table of spaces in the
266           system; they are hashed on the space
267           id */
268   hash_table_t* name_hash;  /*!< hash table based on the space
269           name */
270   UT_LIST_BASE_NODE_T(fil_node_t) LRU;
271           /*!< base node for the LRU list of the
272           most recently used open files with no
273           pending i/o's; if we start an i/o on
274           the file, we first remove it from this
275           list, and return it to the start of
276           the list when the i/o ends;
277           log files and the system tablespace are
278           not put to this list: they are opened
279           after the startup, and kept open until
280           shutdown */
281   UT_LIST_BASE_NODE_T(fil_space_t) unflushed_spaces;
282           /*!< base node for the list of those
283           tablespaces whose files contain
284           unflushed writes; those spaces have
285           at least one file node where
286           modification_counter > flush_counter */
287   ulint    n_open;    /*!< number of files currently open */
288   ulint    max_n_open; /*!< n_open is not allowed to exceed
289           this */
290   ib_int64_t  modification_counter;/*!< when we write to a file we
291           increment this by one */
```

```
292   ulint    max_assigned_id;/*!< maximum space id in the existing
293           tables, or assigned during the time
294           mysqld has been up; at an InnoDB
295           startup we scan the data dictionary
296           and set here the maximum of the
297           space id's of the tables there */
298   ib_int64_t  tablespace_version;
299           /*!< a counter which is incremented for
300           every space object memory creation;
301           every space mem object gets a
302           'timestamp' from this; in DISCARD/
303           IMPORT this is used to check if we
304           should ignore an insert buffer merge
305           request */
306   UT_LIST_BASE_NODE_T(fil_space_t) space_list;
307           /*!< list of all file spaces */
308   ibool    space_id_reuse_warned;
309           /* !< TRUE if fil_space_create()
310           has issued a warning about
311           potential space_id reuse */
312 };
```

# file system init

- `fil/fil0fil.cc:1676, fil_init()`

```
1672 /*****************************************************************//**
1673 Initializes the tablespace memory cache. */
1674 UNIV_INTERN
1675 void
1676 fil_init(
1677 /*=====*/
1678   ulint hash_size,  /*!< in: hash table size */
1679   ulint max_n_open) /*!< in: max number of open files */
1680 {
1681   ut_a(fil_system == NULL);
1682
1683   ut_a(hash_size > 0);
1684   ut_a(max_n_open > 0);
1685
1686   fil_system = static_cast<fil_system_t*>(
1687     mem_zalloc(sizeof(fil_system_t)));
1688
1689   mutex_create(fil_system_mutex_key,
1690         &fil_system->mutex, SYNC_ANY_LATCH);
1691
1692   fil_system->spaces = hash_create(hash_size);
1693   fil_system->name_hash = hash_create(hash_size);
1694
1695   UT_LIST_INIT(fil_system->LRU);
1696
1697   fil_system->max_n_open = max_n_open;
1698 }
```

Protected by mutex

file system hash

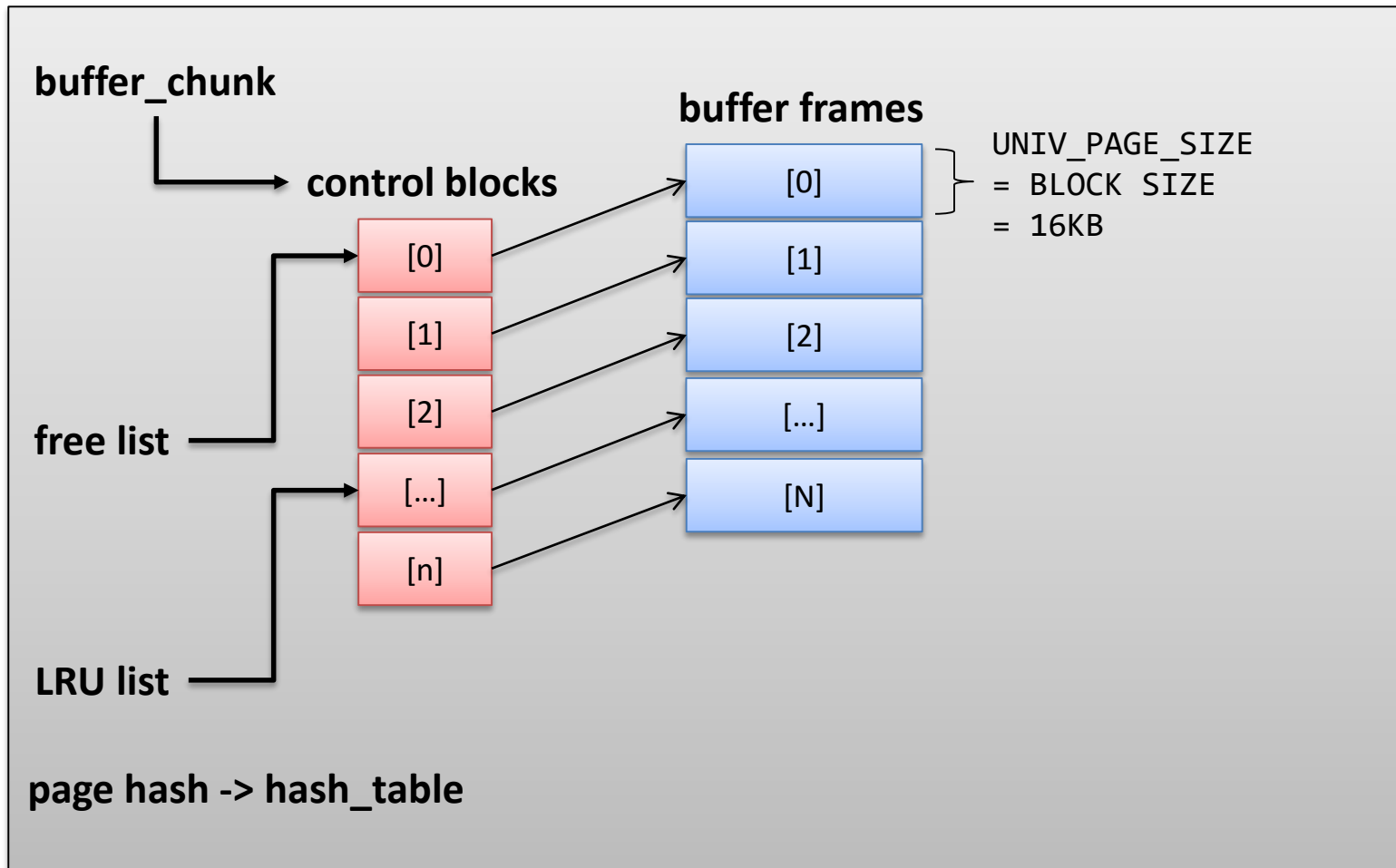Open file LRU

# InnoDB Initialize

- Buffer Pool Init

- `srv0start.cc:1948, buf_pool_init()`

```
1931    fil_init(srv_file_per_table ? 50000 : 5000, srv_max_n_open_files);
1932
1933    double  size;
1934    char    unit;
1935
1936    if (srv_buf_pool_size >= 1024 * 1024 * 1024) {
1937        size = ((double) srv_buf_pool_size) / (1024 * 1024 * 1024);
1938        unit = 'G';
1939    } else {
1940        size = ((double) srv_buf_pool_size) / (1024 * 1024);
1941        unit = 'M';
1942    }
1943
1944    /* Print time to initialize the buffer pool */
1945    ib_logf(IB_LOG_LEVEL_INFO,
1946        "Initializing buffer pool, size = %.1f%c", size, unit);
1947
1948    err = buf_pool_init(srv_buf_pool_size, srv_buf_pool_instances);
```

invoke buffer pool init

# Buffer pool init

**buffer_pool_ptr[]**

**buffer_chunk**

**buffer frames**

**control blocks**

UNIV_PAGE_SIZE
= BLOCK SIZE
= 16KB

| control blocks | buffer frames |
|---|---|
| [0] | [0] |
| [1] | [1] |
| [2] | [2] |
| [...] | [...] |
| [n] | [N] |

**free list**

**LRU list**

**page hash -> hash_table**

# Buffer pool init

- `buf/buf0buf.cc:1384, buf_pool_init()`

```
1384 buf_pool_init(
1385 /*==========*/
1386   ulint total_size, /*!< in: size of the total pool in bytes */
1387   ulint n_instances)  /*!< in: number of instances */
1388 {
1389   ulint   i;
1390   const ulint size  = total_size / n_instances;
1391
1392   ut_ad(n_instances > 0);
1393   ut_ad(n_instances <= MAX_BUFFER_POOLS);
1394   ut_ad(n_instances == srv_buf_pool_instances);
1395
1396   buf_pool_ptr = (buf_pool_t*) mem_zalloc(
1397     n_instances * sizeof *buf_pool_ptr);
1398
1399   for (i = 0; i < n_instances; i++) {
1400     buf_pool_t* ptr = &buf_pool_ptr[i];
1401
1402     if (buf_pool_init_instance(ptr, size, i) != DB_SUCCESS) {
1403
1404       /* Free all the instances created so far. */
1405       buf_pool_free(i);
1406
1407       return(DB_ERROR);
1408     }
1409   }
1410
1411   buf_pool_set_sizes();
1412   buf_LRU_old_ratio_update(100 * 3/ 8, FALSE);
1413
1414   btr_search_sys_create(buf_pool_get_curr_size() / sizeof(void*) / 64);
1415
1416   return(DB_SUCCESS);
1417 }
```

buffer pool init per instance

# Buffer pool init

- buf/buf0buf.cc:1253,
  buf_pool_instance()

```
1248 /*************************************************************//**
1249 Initialize a buffer pool instance.
1250 @return DB_SUCCESS if all goes well. */
1251 UNIV_INTERN
1252 ulint
1253 buf_pool_init_instance(
1254 /*===================*/
1255   buf_pool_t* buf_pool, /*!< in: buffer pool instance */
1256   ulint   buf_pool_size,  /*!< in: size in bytes */
1257   ulint   instance_no)  /*!< in: id of the instance */
1258 {
1259   ulint   i;
1260   buf_chunk_t*  chunk;
1261
1262   /* 1. Initialize general fields
1263   ------------------------------ */
1264   mutex_create(buf_pool_mutex_key,
1265         &buf_pool->mutex, SYNC_BUF_POOL);
1266   mutex_create(buf_pool_zip_mutex_key,
1267         &buf_pool->zip_mutex, SYNC_BUF_BLOCK);
1268
1269   buf_pool_mutex_enter(buf_pool);
1270
1271   if (buf_pool_size > 0) {
1272     buf_pool->n_chunks = 1;
1273
1274     buf_pool->chunks = chunk =
1275       (buf_chunk_t*) mem_zalloc(sizeof *chunk);
1276
1277     UT_LIST_INIT(buf_pool->free);
```

create mutex and enter

buf_pool_zip : use compression

# Buffer pool init

- `buf/buf0buf.cc:1253, buf_pool_instance()`

```
1277    UT_LIST_INIT(buf_pool->free);
1278
1279    if (!buf_chunk_init(buf_pool, chunk, buf_pool_size)) {
1280      mem_free(chunk);
1281      mem_free(buf_pool);
1282
1283      buf_pool_mutex_exit(buf_pool);
1284
1285      return(DB_ERROR);
1286    }
1287
1288    buf_pool->instance_no = instance_no;
1289    buf_pool->old_pool_size = buf_pool_size;
1290    buf_pool->curr_size = chunk->size;
1291    buf_pool->curr_pool_size = buf_pool->curr_size * UNIV_PAGE_SIZE;
```

init buffer chunk - see this later

# Buffer pool init

- `buf/buf0buf.cc:1253,`
  `buf_pool_instance()`

```
1293        /* Number of locks protecting page_hash must be a
1294        power of two */
1295        srv_n_page_hash_locks = static_cast<ulong>(
1296            ut_2_power_up(srv_n_page_hash_locks));
1297        ut_a(srv_n_page_hash_locks != 0);
1298        ut_a(srv_n_page_hash_locks <= MAX_PAGE_HASH_LOCKS);
1299
1300        buf_pool->page_hash = ha_create(2 * buf_pool->curr_size,
1301                srv_n_page_hash_locks,
1302                MEM_HEAP_FOR_PAGE_HASH,
1303                SYNC_BUF_PAGE_HASH);
1304
1305        buf_pool->zip_hash = hash_create(2 * buf_pool->curr_size);
1306
1307        buf_pool->last_printout_time = ut_time();
```
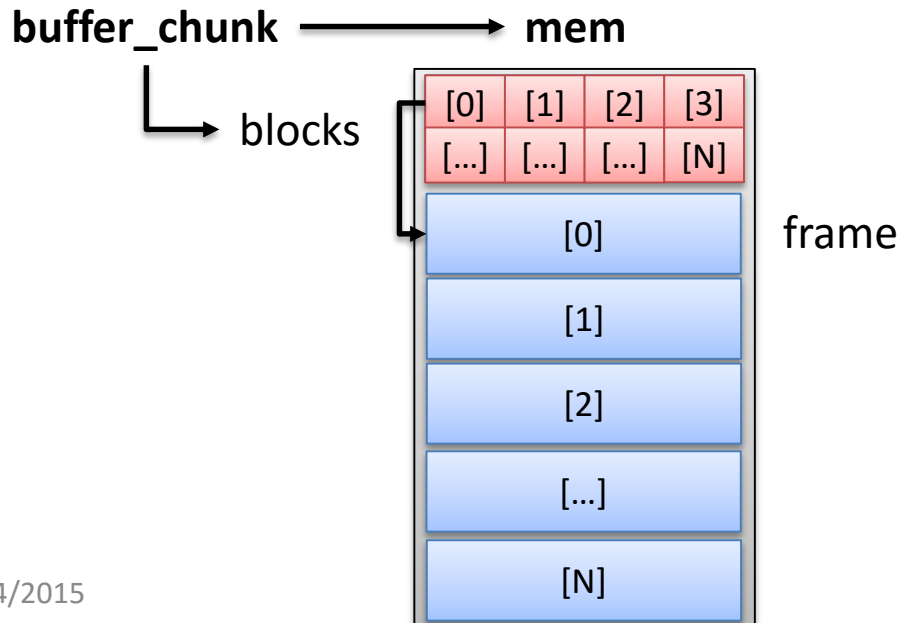
create page hash table

create zip page hash table for compression

# Buffer pool init

- buffer chunk

```
39  /** A chunk of buffers. The buffer pool is allocated in chunks. */
40  struct buf_chunk_t{
41    ulint   mem_size; /*!< allocated size of the chunk */
42    ulint   size;    /*!< size of frames[] and blocks[] */
43    void*   mem;     /*!< pointer to the memory area which
44            was allocated for the frames */
45    buf_block_t*  blocks;   /*!< array of buffer control blocks */
46  };
```
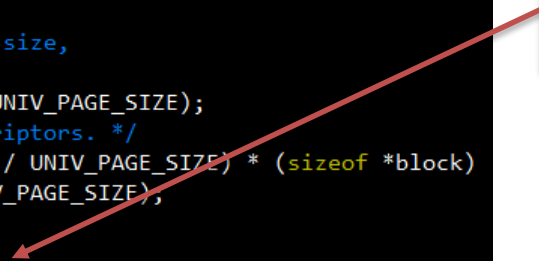
**buffer_chunk** ⟶ **mem**

blocks

| [0] | [1] | [2] | [3] |
|-----|-----|-----|-----|
| [...] | [...] | [...] | [N] |

| [0] |
|-----|

frame

| [1] |
|-----|

| [2] |
|-----|

| [...] |
|-----|

| [N] |
|-----|

# Buffer pool init

- `buf/buf0buf.cc:1041, chunk_init()`

```
1036  /**********************************************************//**
1037  Allocates a chunk of buffer frames.
1038  @return chunk, or NULL on failure */
1039  static
1040  buf_chunk_t*
1041  buf_chunk_init(
1042  /*===========*/
1043    buf_pool_t*  buf_pool, /*!< in: buffer pool instance */
1044    buf_chunk_t*  chunk,    /*!< out: chunk of buffers */
1045    ulint    mem_size) /*!< in: requested size in bytes */
1046  {
1047    buf_block_t*  block;
1048    byte*    frame;
1049    ulint    i;
1050
1051    /* Round down to a multiple of page size,
1052    although it already should be. */
1053    mem_size = ut_2pow_round(mem_size, UNIV_PAGE_SIZE);
1054    /* Reserve space for the block descriptors. */
1055    mem_size += ut_2pow_round((mem_size / UNIV_PAGE_SIZE) * (sizeof *block)
1056            + (UNIV_PAGE_SIZE - 1), UNIV_PAGE_SIZE);
1057
1058    chunk->mem_size = mem_size;
1059    chunk->mem = os_mem_alloc_large(&chunk->mem_size);
1060
1061    if (UNIV_UNLIKELY(chunk->mem == NULL)) {
1062
1063      return(NULL);
1064    }
```

allocate chunk mem (blocks + frames)

# Buffer pool init

- `buf/buf0buf.cc:1041, chunk_init()`

```
1066    /* Allocate the block descriptors from
1067    the start of the memory block. */
1068    chunk->blocks = (buf_block_t*) chunk->mem;
1069
1070    /* Align a pointer to the first frame.  Note that when
1071    os_large_page_size is smaller than UNIV_PAGE_SIZE,
1072    we may allocate one fewer block than requested.  When
1073    it is bigger, we may allocate more blocks than requested. */
1074
1075    frame = (byte*) ut_align(chunk->mem, UNIV_PAGE_SIZE);
1076    chunk->size = chunk->mem_size / UNIV_PAGE_SIZE
1077      - (frame != chunk->mem);
1078
1079    /* Subtract the space needed for block descriptors. */
1080    {
1081      ulint size = chunk->size;
1082
1083      while (frame < (byte*) (chunk->blocks + size)) {
1084        frame += UNIV_PAGE_SIZE;
1085        size--;
1086      }
1087
1088      chunk->size = size;
1089    }
```

allocate control blocks

allocate frame
(page size aligned)

# Buffer pool init

- `buf/buf0buf.cc:1041, chunk_init()`

```
1091    /* Init block structs and assign frames for them. Then we
1092    assign the frames to the first blocks (we already mapped the
1093    memory above). */
1094
1095    block = chunk->blocks;
1096
1097    for (i = chunk->size; i--; ) {
1098
1099        buf_block_init(buf_pool, block, frame);
1100        UNIV_MEM_INVALID(block->frame, UNIV_PAGE_SIZE);
1101
1102        /* Add the block to the free list */
1103        UT_LIST_ADD_LAST(list, buf_pool->free, (&block->page));
1104
1105        ut_d(block->page.in_free_list = TRUE);
1106        ut_ad(buf_pool_from_block(block) == buf_pool);
1107
1108        block++;
1109        frame += UNIV_PAGE_SIZE;
1110    }
1111
1112 #ifdef PFS_GROUP_BUFFER_SYNC
1113    pfs_register_buffer_block(chunk);
1114 #endif
1115    return(chunk);
```

init control block

add all blocks to free list

# Buffer pool init

- `buf/buf0buf.cc:971, buf_block_init()`

```
967 /********************************************************//**
968 Initializes a buffer control block when the buf_pool is created. */
969 static
970 void
971 buf_block_init(
972 /*===========*/
973   buf_pool_t* buf_pool, /*!< in: buffer pool instance */
974   buf_block_t*  block,    /*!< in: pointer to control block */
975   byte*   frame)    /*!< in: pointer to buffer frame */
976 {
977   UNIV_MEM_DESC(frame, UNIV_PAGE_SIZE);
978
979   block->frame = frame;
980
981   block->page.buf_pool_index = buf_pool_index(buf_pool);
982   block->page.state = BUF_BLOCK_NOT_USED;
983   block->page.buf_fix_count = 0;
984   block->page.io_fix = BUF_IO_NONE;
985
986   block->modify_clock = 0;
```

set data frame

page.buf_pool_index : back pointer
page.state : current status
page.buf_fix_count : reference count
page.io_fix : block fix
(Shared/eXclusive)

# Buffer pool init

- `buf/buf0buf.cc:1384, buf_pool_init()`

```
1384 buf_pool_init(
1385 /*==========*/
1386   ulint total_size,  /*!< in: size of the total pool in bytes */
1387   ulint n_instances)  /*!< in: number of instances */
1388 {
1389   ulint   i;
1390   const ulint size = total_size / n_instances;
1391
1392   ut_ad(n_instances > 0);
1393   ut_ad(n_instances <= MAX_BUFFER_POOLS);
1394   ut_ad(n_instances == srv_buf_pool_instances);
1395
1396   buf_pool_ptr = (buf_pool_t*) mem_zalloc(
1397     n_instances * sizeof *buf_pool_ptr);
1398
1399   for (i = 0; i < n_instances; i++) {
1400     buf_pool_t* ptr = &buf_pool_ptr[i];
1401
1402     if (buf_pool_init_instance(ptr, size, i) != DB_SUCCESS) {
1403
1404       /* Free all the instances created so far. */
1405       buf_pool_free(i);
1406
1407       return(DB_ERROR);
1408     }
1409   }
1410
1411   buf_pool_set_sizes();
1412   buf_LRU_old_ratio_update(100 * 3/ 8, FALSE);
1413
1414   btr_search_sys_create(buf_pool_get_curr_size() / sizeof(void*) / 64);
1415
1416   return(DB_SUCCESS);
1417 }
```

set old ratio : 3/8 -> for mid-point insertion

# log init

- `srv/srv0start:1975, log_init()`

```
1974    fsp_init();
1975    log_init();                                                      log init
1976
1977    lock_sys_create(srv_lock_table_size);
1978
1979    /* Create i/o-handler threads: */
1980
1981    for (i = 0; i < srv_n_file_io_threads; ++i) {
1982
1983      n[i] = i;
1984
1985      os_thread_create(io_handler_thread, n + i, thread_ids + i);
1986    }
1987
1988 #ifdef UNIV_LOG_ARCHIVE
1989    if (0 != ut_strcmp(srv_log_group_home_dir, srv_arch_dir)) {
1990      ut_print_timestamp(stderr);
1991      fprintf(stderr, " InnoDB: Error: you must set the log group home dir in my.cnf\n");
1992      ut_print_timestamp(stderr);
1993      fprintf(stderr, " InnoDB: the same as log arch dir.\n");
1994
1995      return(DB_ERROR);
1996    }
1997 #endif /* UNIV_LOG_ARCHIVE */
1998
```

# log init

- `log/log0log.cc:838, log_init()`

```
835 Initializes the log. */
836 UNIV_INTERN
837 void
838 log_init(void)
839 /*==========*/
840 {
841   log_sys = static_cast<log_t*>(mem_alloc(sizeof(log_t)));
842
843   mutex_create(log_sys_mutex_key, &log_sys->mutex, SYNC_LOG);
844
845   mutex_create(log_flush_order_mutex_key,
846           &log_sys->log_flush_order_mutex,
847           SYNC_LOG_FLUSH_ORDER);
848
849   mutex_enter(&(log_sys->mutex));
850
851   /* Start the lsn from one log block from zero: this way every
852   log record has a start lsn != zero, a fact which we will use */
853
854   log_sys->lsn = LOG_START_LSN;
855
856   ut_a(LOG_BUFFER_SIZE >= 16 * OS_FILE_LOG_BLOCK_SIZE);
857   ut_a(LOG_BUFFER_SIZE >= 4 * UNIV_PAGE_SIZE);
858
859   log_sys->buf_ptr = static_cast<byte*>(
860     mem_zalloc(LOG_BUFFER_SIZE + OS_FILE_LOG_BLOCK_SIZE));
861
862   log_sys->buf = static_cast<byte*>(
863     ut_align(log_sys->buf_ptr, OS_FILE_LOG_BLOCK_SIZE));
864
865   log_sys->buf_size = LOG_BUFFER_SIZE;
866   log_sys->is_extending = false;
```

create log sys

mutex for protection and synchronization

start_lsn != 0

buf_ptr : mem buffer
**buf : page aligned buffer for log**

# log init

- `log/log0log.cc:838, log_init()`

```
879   log_sys->buf_next_to_write = 0;
880
881   log_sys->write_lsn = 0;
882   log_sys->current_flush_lsn = 0;
883   log_sys->flushed_to_disk_lsn = 0;
884
885   log_sys->written_to_some_lsn = log_sys->lsn;
886   log_sys->written_to_all_lsn = log_sys->lsn;
887
888   log_sys->n_pending_writes = 0;
889
890   log_sys->no_flush_event = os_event_create();
891
892   os_event_set(log_sys->no_flush_event);
893
894   log_sys->one_flushed_event = os_event_create();
895
896   os_event_set(log_sys->one_flushed_event);
897
898   /*----------------------------*/
899
900   log_sys->next_checkpoint_no = 0;
901   log_sys->last_checkpoint_lsn = log_sys->lsn;
902   log_sys->n_pending_checkpoint_writes = 0;
903
904
905   rw_lock_create(checkpoint_lock_key, &log_sys->checkpoint_lock,
906          SYNC_NO_ORDER_CHECK);
907
908   log_sys->checkpoint_buf_ptr = static_cast<byte*>(
909     mem_zalloc(2 * OS_FILE_LOG_BLOCK_SIZE));
910
911   log_sys->checkpoint_buf = static_cast<byte*>(
912     ut_align(log_sys->checkpoint_buf_ptr, OS_FILE_LOG_BLOCK_SIZE));
```

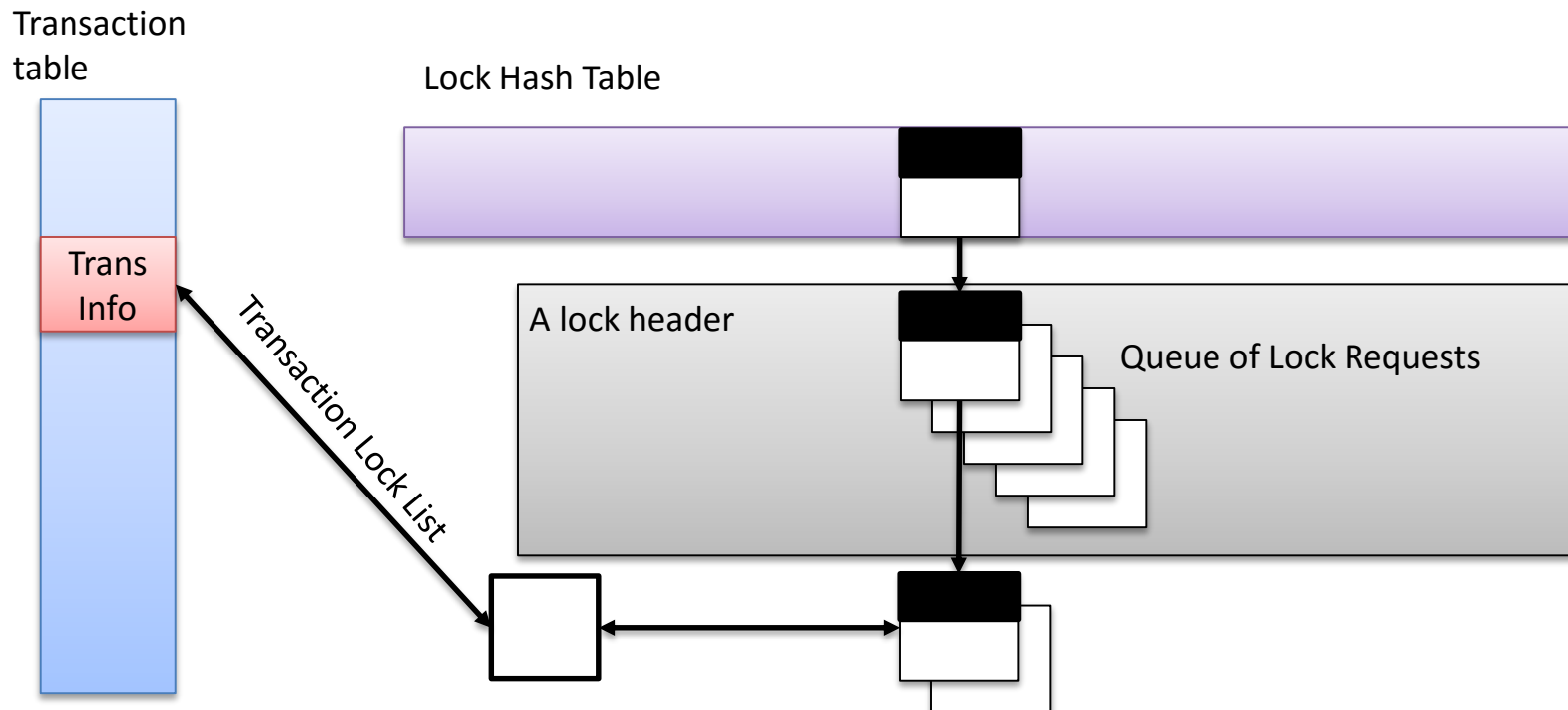init log sys variables

# Contents

- Before we begin : tools
  - message logging and entry point
  - Coding conventions : the very first

- InnoDB Initialize
  - Init. in-memory objects and check configurations
  - aio system, file, buffer pool, log system

- Create objects
  - lock-system, io-handler, data files and log files

- Open
  - log and system table space
  - transaction system

- Recovery if necessary

- Create double write system if necessary

# CREATE SERVER OBJECTS

# Locking system

- Locking (and transaction model)
  - When a transaction updates a row in a table, or locks it with SELECT FOR UPDATE, InnoDB establishes a list or queue of locks on that row.
  - Similarly, InnoDB maintains a list of locks on a table for table-level locks transactions hold.

Transaction table

Lock Hash Table

Trans Info

Transaction Lock List

A lock header

Queue of Lock Requests

# Locking system

- ## Lock Compatibility
  - S : shared, X: eXclusive

| | Granted Mode | | | | | | |
|---|---|---|---|---|---|---|---|
| Requested Mode | None | IS | IX | S | SIX | U | X |
| IS | + | + | + | + | + | - | - |
| IX | + | + | + | - | - | - | - |
| S | + | + | - | + | - | - | - |
| SIX | + | + | - | - | - | - | - |
| U | + | - | - | + | - | - | - |
| X | + | - | - | - | - | - | - |

# Locking system

- `srv/srv0start:1977, lock_sys_create()`

```
1974   fsp_init();
1975   log_init();
1976
1977   lock_sys_create(srv_lock_table_size);        ← lock system create
1978
1979   /* Create i/o-handler threads: */
1980
1981   for (i = 0; i < srv_n_file_io_threads; ++i) {
1982
1983     n[i] = i;
1984
1985     os_thread_create(io_handler_thread, n + i, thread_ids + i);
1986   }
1987
1988 #ifdef UNIV_LOG_ARCHIVE
1989   if (0 != ut_strcmp(srv_log_group_home_dir, srv_arch_dir)) {
1990     ut_print_timestamp(stderr);
1991     fprintf(stderr, " InnoDB: Error: you must set the log group home dir in my.cnf\n");
1992     ut_print_timestamp(stderr);
1993     fprintf(stderr, " InnoDB: the same as log arch dir.\n");
1994
1995     return(DB_ERROR);
1996   }
1997 #endif /* UNIV_LOG_ARCHIVE */
1998
```

# Locking system

- lock table size
  - in `srv_normalize_init_values()`
  - 5 * # of pages in the buffer pool

```
/***********************************************************************//**
Normalizes init parameter values to use units we use inside InnoDB. */
static
void
srv_normalize_init_values(void)
/*===========================*/
{
  ulint n;
  ulint i;

  n = srv_n_data_files;

  for (i = 0; i < n; i++) {
    srv_data_file_sizes[i] = srv_data_file_sizes[i]
      * ((1024 * 1024) / UNIV_PAGE_SIZE);
  }

  srv_last_file_size_max = srv_last_file_size_max
    * ((1024 * 1024) / UNIV_PAGE_SIZE);

  srv_log_file_size = srv_log_file_size / UNIV_PAGE_SIZE;

  srv_log_buffer_size = srv_log_buffer_size / UNIV_PAGE_SIZE;

  srv_lock_table_size = 5 * (srv_buf_pool_size / UNIV_PAGE_SIZE);
}
```

lock table size

# Lock system

- `lock/lock0lock:589, lock_sys_create()`

```
/*******************************************************************//**
Creates the lock system at database start. */
UNIV_INTERN
void
lock_sys_create(
/*============*/
  ulint n_cells)  /*!< in: number of slots in lock hash table */
{
  ulint lock_sys_sz;

  lock_sys_sz = sizeof(*lock_sys)
    + OS_THREAD_MAX_N * sizeof(srv_slot_t);

  lock_sys = static_cast<lock_sys_t*>(mem_zalloc(lock_sys_sz));

  lock_stack = static_cast<lock_stack_t*>(
    mem_zalloc(sizeof(*lock_stack) * LOCK_STACK_SIZE));

  void* ptr = &lock_sys[1];

  lock_sys->waiting_threads = static_cast<srv_slot_t*>(ptr);

  lock_sys->last_slot = lock_sys->waiting_threads;

  mutex_create(lock_sys_mutex_key, &lock_sys->mutex, SYNC_LOCK_SYS);

  mutex_create(lock_sys_wait_mutex_key,
        &lock_sys->wait_mutex, SYNC_LOCK_WAIT_SYS);

  lock_sys->timeout_event = os_event_create();

  lock_sys->rec_hash = hash_create(n_cells);

  if (!srv_read_only_mode) {
    lock_latest_err_file = os_file_create_tmpfile();
    ut_a(lock_latest_err_file);
  }
}
```

create record hash

# Create io threads

- `srv/srv0start:1985`

```
1974    fsp_init();
1975    log_init();
1976
1977    lock_sys_create(srv_lock_table_size);
1978
1979    /* Create i/o-handler threads: */
1980
1981    for (i = 0; i < srv_n_file_io_threads; ++i) {
1982
1983      n[i] = i;
1984
1985      os_thread_create(io_handler_thread, n + i, thread_ids + i);
1986    }
1987
1988 #ifdef UNIV_LOG_ARCHIVE
1989    if (0 != ut_strcmp(srv_log_group_home_dir, srv_arch_dir)) {
1990      ut_print_timestamp(stderr);
1991      fprintf(stderr, " InnoDB: Error: you must set the log group home dir in my.cnf\n");
1992      ut_print_timestamp(stderr);
1993      fprintf(stderr, " InnoDB: the same as log arch dir.\n");
1994
1995      return(DB_ERROR);
1996    }
1997 #endif /* UNIV_LOG_ARCHIVE */
1998
```

create io threads

# Recovery system

- Recovery system
  - It mainly focuses on write-ahead-log and double write
  - If the system was not clean shutdown, then we have to do recovery

# Recovery system

- srv/srv0start:2056, recv_sys_create()

```
2056    recv_sys_create();
2057    recv_sys_init(buf_pool_get_curr_size());
2058
2059    err = open_or_create_data_files(&create_new_db,
2060 #ifdef UNIV_LOG_ARCHIVE
2061            &min_arch_log_no, &max_arch_log_no,
2062 #endif /* UNIV_LOG_ARCHIVE */
2063            &min_flushed_lsn, &max_flushed_lsn,
2064            &sum_of_new_sizes);
2065    if (err == DB_FAIL) {
2066
2067        ib_logf(IB_LOG_LEVEL_ERROR,
2068            "The system tablespace must be writable!");
2069
2070        return(DB_ERROR);
2071
2072    } else if (err != DB_SUCCESS) {
2073
2074        ib_logf(IB_LOG_LEVEL_ERROR,
2075            "Could not open or create the system tablespace. If "
2076            "you tried to add new data files to the system "
2077            "tablespace, and it failed here, you should now "
2078            "edit innodb_data_file_path in my.cnf back to what "
2079            "it was, and remove the new ibdata files InnoDB "
2080            "created in this failed attempt. InnoDB only wrote "
2081            "those files full of zeros, but did not yet use "
2082            "them in any way. But be careful: do not remove "
2083            "old data files which contain your precious data!");
2084
2085        return(err);
2086    }
```

create recovery system

# Recovery system

- `log/log0recv.cc:196, recv_sys_create()`
  - just make `recv_sys`

```
192 /*****************************************************//**
193 Creates the recovery system. */
194 UNIV_INTERN
195 void
196 recv_sys_create(void)
197 /*=================*/
198 {
199   if (recv_sys != NULL) {
200
201     return;
202   }
203
204   recv_sys = static_cast<recv_sys_t*>(mem_zalloc(sizeof(*recv_sys)));
205
206   mutex_create(recv_sys_mutex_key, &recv_sys->mutex, SYNC_RECV);
207
208 #ifndef UNIV_HOTBACKUP
209   mutex_create(recv_writer_mutex_key, &recv_sys->writer_mutex,
210         SYNC_LEVEL_VARYING);
211 #endif /* !UNIV_HOTBACKUP */
212
213   recv_sys->heap = NULL;
214   recv_sys->addr_hash = NULL;
215 }
```

# Recovery system

- `srv/srv0start:2057, recv_sys_init()`

```
2056    recv_sys_create();
2057    recv_sys_init(buf_pool_get_curr_size());
2058
2059    err = open_or_create_data_files(&create_new_db,
2060 #ifdef UNIV_LOG_ARCHIVE
2061            &min_arch_log_no, &max_arch_log_no,
2062 #endif /* UNIV_LOG_ARCHIVE */
2063            &min_flushed_lsn, &max_flushed_lsn,
2064            &sum_of_new_sizes);
2065    if (err == DB_FAIL) {
2066
2067        ib_logf(IB_LOG_LEVEL_ERROR,
2068            "The system tablespace must be writable!");
2069
2070        return(DB_ERROR);
2071
2072    } else if (err != DB_SUCCESS) {
2073
2074        ib_logf(IB_LOG_LEVEL_ERROR,
2075            "Could not open or create the system tablespace. If "
2076            "you tried to add new data files to the system "
2077            "tablespace, and it failed here, you should now "
2078            "edit innodb_data_file_path in my.cnf back to what "
2079            "it was, and remove the new ibdata files InnoDB "
2080            "created in this failed attempt. InnoDB only wrote "
2081            "those files full of zeros, but did not yet use "
2082            "them in any way. But be careful: do not remove "
2083            "old data files which contain your precious data!");
2084
2085        return(err);
2086    }
```

init recovery system

# Recovery system

- Definition recovery system
  - `include/log0recv.h:384,recv_sys_t`

```
383 /** Recovery system data structure */
384 struct recv_sys_t{
385 #ifndef UNIV_HOTBACKUP
386   ib_mutex_t     mutex;  /*!< mutex protecting the fields apply_log_recs,
387       n_addrs, and the state field in each recv_addr
388       struct */
389   ib_mutex_t     writer_mutex;/*!< mutex coordinating
390       flushing between recv_writer_thread and
391       the recovery thread. */
392 #endif /* !UNIV_HOTBACKUP */
393   ibool   apply_log_recs;
394       /*!< this is TRUE when log rec application to
395       pages is allowed; this flag tells the
396       i/o-handler if it should do log record
397       application */
398   ibool   apply_batch_on;
399       /*!< this is TRUE when a log rec application
400       batch is running */
401   lsn_t   lsn;  /*!< log sequence number */
402   ulint   last_log_buf_size;
403       /*!< size of the log buffer when the database
404       last time wrote to the log */
405   byte*   last_block;
406       /*!< possible incomplete last recovered log
407       block */
408   byte*   last_block_buf_start;
409       /*!< the nonaligned start address of the
410       preceding buffer */
```

# Recovery system

- Definition recovery system
  - `include/log0recv.h:384,recv_sys_t`

```
411    byte*   buf;   /*!< buffer for parsing log records */
412    ulint   len;   /*!< amount of data in buf */
413    lsn_t   parse_start_lsn;
414          /*!< this is the lsn from which we were able to
415          start parsing log records and adding them to
416          the hash table; zero if a suitable
417          start point not found yet */
418    lsn_t   scanned_lsn;
419          /*!< the log data has been scanned up to this
420          lsn */
421    ulint   scanned_checkpoint_no;
422          /*!< the log data has been scanned up to this
423          checkpoint number (lowest 4 bytes) */
424    ulint   recovered_offset;
425          /*!< start offset of non-parsed log records in
426          buf */
427    lsn_t   recovered_lsn;
428          /*!< the log records have been parsed up to
429          this lsn */
430    lsn_t   limit_lsn;/*!< recovery should be made at most
431          up to this lsn */
432    ibool   found_corrupt_log;
433          /*!< this is set to TRUE if we during log
434          scan find a corrupt log block, or a corrupt
435          log record, or there is a log parsing
436          buffer overflow */
```

# Recovery system

- Definition recovery system
  - `include/log0recv.h:384,recv_sys_t`

```
437 #ifdef UNIV_LOG_ARCHIVE
438    log_group_t*  archive_group;
439         /*!< in archive recovery: the log group whose
440         archive is read */
441 #endif /* !UNIV_LOG_ARCHIVE */
442    mem_heap_t* heap; /*!< memory heap of log records and file
443         addresses*/
444    hash_table_t* addr_hash;/*!< hash table of file addresses of pages */
445    ulint   n_addrs;/*!< number of not processed hashed file
446         addresses in the hash table */
447
448    recv_dblwr_t  dblwr;
449 };
```

# Recovery system

- `log/log0recv.cc:380, recv_sys_init()`

```
376 /*************************************************************
377 Inits the recovery system for a recovery operation. */
378 UNIV_INTERN
379 void
380 recv_sys_init(
381 /*==========*/
382   ulint available_memory) /*!< in: available memory in bytes */
383 {
384   if (recv_sys->heap != NULL) {
385
386     return;
387   }
388
389 #ifndef UNIV_HOTBACKUP
390   /* Initialize red-black tree for fast insertions into the
391   flush_list during recovery process.
392   As this initialization is done while holding the buffer pool
393   mutex we perform it before acquiring recv_sys->mutex. */
394   buf_flush_init_flush_rbt();                          make flush list: see the details later
395
396   mutex_enter(&(recv_sys->mutex));
397
398   recv_sys->heap = mem_heap_create_typed(256,
399         MEM_HEAP_FOR_RECV_SYS);
400 #else /* !UNIV_HOTBACKUP */
401   recv_sys->heap = mem_heap_create(256);
402   recv_is_from_backup = TRUE;
403 #endif /* !UNIV_HOTBACKUP */
```

# Recovery system

- `log/log0recv.cc:380, recv_sys_init()`

```
404
405    /* Set appropriate value of recv_n_pool_free_frames. */
406    if (buf_pool_get_curr_size() >= (10 * 1024 * 1024)) {
407      /* Buffer pool of size greater than 10 MB. */
408      recv_n_pool_free_frames = 512;
409    }
410
411    recv_sys->buf = static_cast<byte*>(ut_malloc(RECV_PARSING_BUF_SIZE));
412    recv_sys->len = 0;
413    recv_sys->recovered_offset = 0;
414
415    recv_sys->addr_hash = hash_create(available_memory / 512);
416    recv_sys->n_addrs = 0;
417
418    recv_sys->apply_log_recs = FALSE;
419    recv_sys->apply_batch_on = FALSE;
420
421    recv_sys->last_block_buf_start = static_cast<byte*>(
422      mem_alloc(2 * OS_FILE_LOG_BLOCK_SIZE));
423
424    recv_sys->last_block = static_cast<byte*>(ut_align(
425      recv_sys->last_block_buf_start, OS_FILE_LOG_BLOCK_SIZE));
426
427    recv_sys->found_corrupt_log = FALSE;
428
429    recv_max_page_lsn = 0;
430
431    /* Call the constructor for recv_sys_t::dblwr member */
432    new (&recv_sys->dblwr) recv_dblwr_t();
433
434    mutex_exit(&(recv_sys->mutex));
435 }
```

make recovery system buffer pool for log scan

Init default variables and make buffers

create double write system

67

# OPEN DATA FILES AND LOG FILE

# check data files

- `srv/srv0start.cc:2059`

```
2059    err = open_or_create_data_files(&create_new_db,
2060 #ifdef UNIV_LOG_ARCHIVE
2061          &min_arch_log_no, &max_arch_log_no,
2062 #endif /* UNIV_LOG_ARCHIVE */
2063          &min_flushed_lsn, &max_flushed_lsn,
2064          &sum_of_new_sizes);
2065   if (err == DB_FAIL) {
2066
2067     ib_logf(IB_LOG_LEVEL_ERROR,
2068       "The system tablespace must be writable!");
2069
2070     return(DB_ERROR);
2071
2072   } else if (err != DB_SUCCESS) {
2073
2074     ib_logf(IB_LOG_LEVEL_ERROR,
2075       "Could not open or create the system tablespace. If "
2076       "you tried to add new data files to the system "
2077       "tablespace, and it failed here, you should now "
2078       "edit innodb_data_file_path in my.cnf back to what "
2079       "it was, and remove the new ibdata files InnoDB "
2080       "created in this failed attempt. InnoDB only wrote "
2081       "those files full of zeros, but did not yet use "
2082       "them in any way. But be careful: do not remove "
2083       "old data files which contain your precious data!");
2084
2085     return(err);
2086   }
```

Open data files

open data files error handling

# check data files

- srv/srv0start.cc:759, open_or_create_data_files()
  - Skip

```
989        /* This is the earliest location where we can load
990        the double write buffer. */
991        if (i == 0) {
992          buf_dblwr_init_or_load_pages(
993            files[i], srv_data_file_names[i], true);
994        }
```

init or load double write pages

# check log files

- srv/srv0start.cc

```
2119    for (i = 0; i < SRV_N_LOG_FILES_MAX; i++) {
2120        os_offset_t size;
2121        os_file_stat_t  stat_info;
2122
2123        sprintf(logfilename + dirnamelen,
2124            "ib_logfile%u", i);
2125
2126        err = os_file_get_status(
2127            logfilename, &stat_info, false);
2128
2129        if (err == DB_NOT_FOUND) {
2130            if (i == 0) {
2131                if (max_flushed_lsn
2132                    != min_flushed_lsn) {
2133                    ib_logf(IB_LOG_LEVEL_ERROR,
2134                        "Cannot create"
2135                        " log files because"
2136                        " data files are"
2137                        " corrupt or"
2138                        " not in sync"
2139                        " with each other");
2140                    return(DB_ERROR);
2141                }
```

get log file status

if log file not found, then create new one

# check log files

- `srv/srv0start.cc:2189`



get log file status

```
2185        if (!srv_file_check_mode(logfilename)) {
2186            return(DB_ERROR);
2187        }
2188
2189        err = open_log_file(&files[i], logfilename, &size);
2190
2191        if (err != DB_SUCCESS) {
2192            return(err);
2193        }
2194
2195        ut_a(size != (os_offset_t) -1);
2196
2197        if (size & ((1 << UNIV_PAGE_SIZE_SHIFT) - 1)) {
2198            ib_logf(IB_LOG_LEVEL_ERROR,
2199                "Log file %s size "
2200                UINT64PF " is not a multiple of"
2201                " innodb_page_size",
2202                logfilename, size);
2203            return(DB_ERROR);
2204        }
```
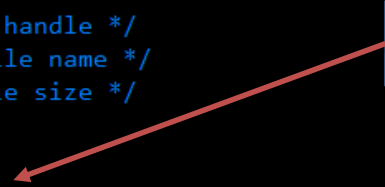
# check log files

- `srv/srv0start.cc:3219, open_log_file()`

```
726  /*******************************************************************/ /**
727  Opens a log file.
728  @return DB_SUCCESS or error code */
729  static __attribute__((nonnull, warn_unused_result))
730  dberr_t
731  open_log_file(
732  /*==========*/
733    os_file_t*   file,  /*!< out: file handle */
734    const char*  name,  /*!< in: log file name */
735    os_offset_t*  size)  /*!< out: file size */
736  {
737    ibool ret;
738
739    *file = os_file_create(innodb_file_log_key, name,
740              OS_FILE_OPEN, OS_FILE_AIO,
741              OS_LOG_FILE, &ret);
742    if (!ret) {
743      ib_logf(IB_LOG_LEVEL_ERROR, "Unable to open '%s'", name);
744      return(DB_ERROR);
745    }
746
747    *size = os_file_get_size(*file);
748
749    ret = os_file_close(*file);
750    ut_a(ret);
751    return(DB_SUCCESS);
752  }
```

create mode = OS_FILE_OPEN

# check log files

- `srv/srv0start.cc:2230`

```
2226      /* Create the in-memory file space objects. */
2227
2228      sprintf(logfilename + dirnamelen, "ib_logfile%u", 0);
2229
2230      fil_space_create(logfilename,
2231          SRV_LOG_SPACE_FIRST_ID,
2232          fsp_flags_set_page_size(0, UNIV_PAGE_SIZE),
2233          FIL_LOG);
2234
2235      ut_a(fil_validate());
2236
2237      /* srv_log_file_size is measured in pages; if page size is 16KB,
2238      then we have a limit of 64TB on 32 bit systems */
2239      ut_a(srv_log_file_size <= ULINT_MAX);
2240
2241      for (unsigned j = 0; j < i; j++) {
2242          sprintf(logfilename + dirnamelen, "ib_logfile%u", j);
2243
2244          if (!fil_node_create(logfilename,
2245                  (ulint) srv_log_file_size,
2246                  SRV_LOG_SPACE_FIRST_ID, FALSE)) {
2247              return(DB_ERROR);
2248          }
```

file system create for log files

file node create for log files

# Open log and system table space

- `srv/srv0start.cc:2267,`
  `fil_open_log_and_system_tablespace_files()`

```
2262 files_checked:
2263    /* Open all log files and data files in the system
2264    tablespace: we keep them open until database
2265    shutdown */
2266
2267    fil_open_log_and_system_tablespace_files();
2268
2269    err = srv_undo_tablespaces_init(
2270        create_new_db,
2271        srv_undo_tablespaces,
2272        &srv_undo_tablespaces_open);
2273
2274    /* If the force recovery is set very high then we carry on regardless
2275    of all errors. Basically this is fingers crossed mode. */
2276
2277    if (err != DB_SUCCESS
2278        && srv_force_recovery < SRV_FORCE_NO_UNDO_LOG_SCAN) {
2279
2280        return(err);
2281    }
```

file node create for log files

# Open log and system table space

- `fil/fil0fil.cc:1708,`
  `fil_open_log_and_system_tablespace_files()`

```
1706 UNIV_INTERN
1707 void
1708 fil_open_log_and_system_tablespace_files(void)
1709 /*==========================================*/
1710 {
1711   fil_space_t*   space;
1712
1713   mutex_enter(&fil_system->mutex);
1714
1715   for (space = UT_LIST_GET_FIRST(fil_system->space_list);
1716        space != NULL;
1717        space = UT_LIST_GET_NEXT(space_list, space)) {
1718
1719     fil_node_t* node;
1720
1721     if (fil_space_belongs_in_lru(space)) {
1722
1723       continue;
1724     }
```

# Open log and system table space

- `fil/fil0fil.cc:1708,`
  `fil_open_log_and_system_tablespace_files()`

```
1726      for (node = UT_LIST_GET_FIRST(space->chain);
1727          node != NULL;
1728          node = UT_LIST_GET_NEXT(chain, node)) {
1729
1730        if (!node->open) {
1731          if (!fil_node_open_file(node, fil_system,
1732              space)) {
1733            /* This func is called during server's
1734            startup. If some file of log or system
1735            tablespace is missing, the server
1736            can't start successfully. So we should
1737            assert for it. */
1738            ut_a(0);
1739          }
1740        }
1764      }
1765    }
1766
1767  mutex_exit(&fil_system->mutex);
1768 }
```

file node open

# INIT UNDO TABLESPACE

# Undo table space init

- `srv/srv0start.cc:2269,`
  `srv_undo_tablespace_init()`

```
2262 files_checked:
2263    /* Open all log files and data files in the system
2264    tablespace: we keep them open until database
2265    shutdown */
2266
2267    fil_open_log_and_system_tablespace_files();
2268
2269    err = srv_undo_tablespaces_init(
2270        create_new_db,
2271        srv_undo_tablespaces,
2272        &srv_undo_tablespaces_open);
2273
2274    /* If the force recovery is set very high then we carry on regardless
2275    of all errors. Basically this is fingers crossed mode. */
2276
2277    if (err != DB_SUCCESS
2278        && srv_force_recovery < SRV_FORCE_NO_UNDO_LOG_SCAN) {
2279
2280        return(err);
2281    }
```

undo table space init

# Undo table space init

- srv/srv0start.cc:1272,
  srv_undo_tablespace_init()

```
1270 static
1271 dberr_t
1272 srv_undo_tablespaces_init(
1273 /*======================*/
1274   ibool    create_new_db,     /*!< in: TRUE if new db being
1275              created */
1276   const ulint n_conf_tablespaces, /*!< in: configured undo
1277              tablespaces */
1278   ulint*   n_opened)   /*!< out: number of UNDO
1279              tablespaces successfully
1280              discovered and opened */
1281 {
1282   ulint   i;
1283   dberr_t   err = DB_SUCCESS;
1284   ulint   prev_space_id = 0;
1285   ulint   n_undo_tablespaces;
1286   ulint   undo_tablespace_ids[TRX_SYS_N_RSEGS + 1];
1287
1288   *n_opened = 0;
1289
1290   ut_a(n_conf_tablespaces <= TRX_SYS_N_RSEGS);
1291
1292   memset(undo_tablespace_ids, 0x0, sizeof(undo_tablespace_ids));
```

# Undo table space init

- srv/srv0start.cc:1272,
  srv_undo_tablespace_init()

```
1294    /* Create the undo spaces only if we are creating a new
1295    instance. We don't allow creating of new undo tablespaces
1296    in an existing instance (yet).  This restriction exists because
1297    we check in several places for SYSTEM tablespaces to be less than
1298    the min of user defined tablespace ids. Once we implement saving
1299    the location of the undo tablespaces and their space ids this
1300    restriction will/should be lifted. */
1301
1302    for (i = 0; create_new_db && i < n_conf_tablespaces; ++i) {
1303      char  name[OS_FILE_MAX_PATH];
1304
1305      ut_snprintf(
1306        name, sizeof(name),
1307        "%s%cundo%03lu",
1308        srv_undo_dir, SRV_PATH_SEPARATOR, i + 1);
1309
1310      /* Undo space ids start from 1. */
1311      err = srv_undo_tablespace_create(
1312        name, SRV_UNDO_TABLESPACE_SIZE_IN_PAGES);
1313
1314      if (err != DB_SUCCESS) {
1315
1316        ib_logf(IB_LOG_LEVEL_ERROR,
1317          "Could not create undo tablespace '%s'.",
1318          name);
1319
1320        return(err);
1321      }
1322    }
```

create new db case, skip

# Undo table space init

- `srv/srv0start.cc:1272, srv_undo_tablespace_init()`

```
1347    for (i = 0; i < n_undo_tablespaces; ++i) {
1348      char   name[OS_FILE_MAX_PATH];
1349
1363      /* Undo space ids start from 1. */
1364
1365      err = srv_undo_tablespace_open(name, undo_tablespace_ids[i]);
1366
1367      if (err != DB_SUCCESS) {
1368
1369        ib_logf(IB_LOG_LEVEL_ERROR,
1370          "Unable to open undo tablespace '%s'.", name);
1371
1372        return(err);
1373      }
1374
1375      prev_space_id = undo_tablespace_ids[i];
1376
1377      ++*n_opened;
1378    }
```

open undo tablespace

# START TRANSACTION SYSTEM

# Create transaction system

- `srv/srv0start.cc:2291, trx_sys_create()`

```
2282
2283    /* Initialize objects used by dict stats gathering thread, which
2284    can also be used by recovery if it tries to drop some table */
2285    if (!srv_read_only_mode) {
2286      dict_stats_thread_init();
2287    }
2288
2289    trx_sys_file_format_init();
2290
2291    trx_sys_create();
2292
2293    if (create_new_db) {
2294
2295      ut_a(!srv_read_only_mode);
2296
2297      mtr_start(&mtr);
2298
2299      fsp_header_init(0, sum_of_new_sizes, &mtr);
2300
2301      mtr_commit(&mtr);
```

create transaction system

# Create transaction system

- `trx/trx0sys.cc:587, trx_sys_create()`

```
583 /*********************************************************//**
584 Creates the trx_sys instance and initializes ib_bh and mutex. */
585 UNIV_INTERN
586 void
587 trx_sys_create(void)
588 /*===============*/
589 {
590   ut_ad(trx_sys == NULL);
591
592   trx_sys = static_cast<trx_sys_t*>(mem_zalloc(sizeof(*trx_sys)));
593
594   mutex_create(trx_sys_mutex_key, &trx_sys->mutex, SYNC_TRX_SYS);
595 }
596
```

# Do recovery

- `srv/srv0start.cc:2404,`
  `recv_recovery_from_checkpoint_start()`

```
2401    /* We always try to do a recovery, even if the database had
2402    been shut down normally: this is the normal startup path */
2403
2404    err = recv_recovery_from_checkpoint_start(
2405      LOG_CHECKPOINT, LSN_MAX,
2406      min_flushed_lsn, max_flushed_lsn);
2407
2408    if (err != DB_SUCCESS) {
2409
2410      return(DB_ERROR);
2411    }
2412
```

recovery from checkpoint:
see details later

# Initialize transaction system

- `srv/srv0start.cc:2425,`
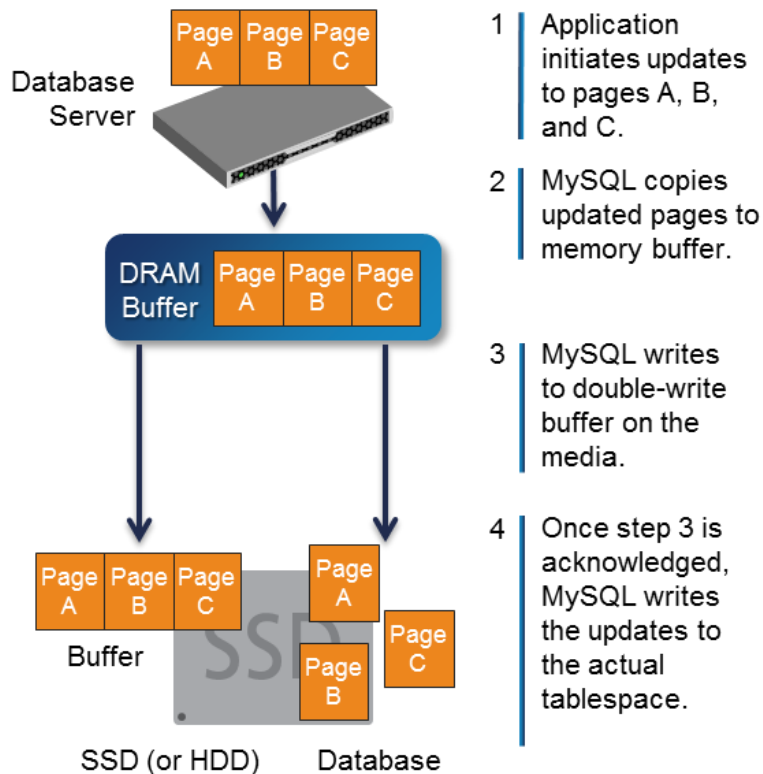  `trx_sys_init_at_db_start()`

see details later

```
2425    ib_bh = trx_sys_init_at_db_start();
2426    n_recovered_trx = UT_LIST_GET_LEN(trx_sys->rw_trx_list);
2427
2428    /* The purge system needs to create the purge view and
2429    therefore requires that the trx_sys is inited. */
2430
2431    trx_purge_sys_create(srv_n_purge_threads, ib_bh);
2432
2433    /* recv_recovery_from_checkpoint_finish needs trx lists which
2434    are initialized in trx_sys_init_at_db_start(). */
2435
2436    recv_recovery_from_checkpoint_finish();
2437
```
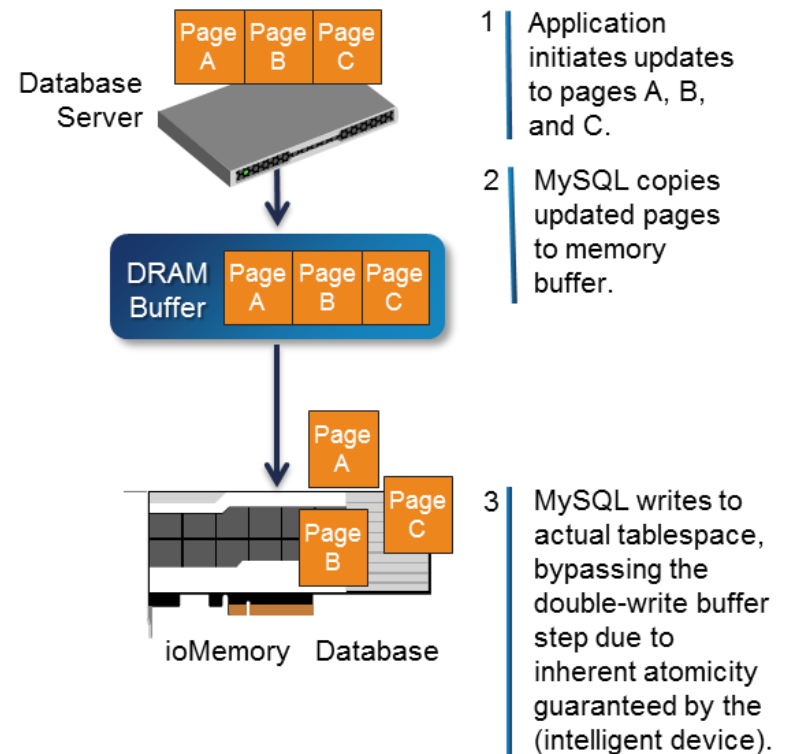
# CREATE DOUBLE WRITE BUFFER

# Double write buffer



http://www.fusionio.com/load/-media-/29ek9p/imagesLibrary/03_atomic_writes.png

# double write create

- `srv/srv0start.cc:2605, buf_dblwr_create()`

```
2601
2602    if (buf_dblwr == NULL) {
2603      /* Create the doublewrite buffer to a new tablespace */
2604
2605      buf_dblwr_create();
2606    }
2607
2608    /* Here the double write buffer has already been created and so
2609    any new rollback segments will be allocated after the double
2610    write buffer. The default segment should already exist.
2611    We create the new segments only if it's a new database or
2612    the database was shutdown cleanly. */
2613
2614    /* Note: When creating the extra rollback segments during an upgrade
2615    we violate the latching order, even if the change buffer is empty.
2616    We make an exception in sync0sync.cc and check srv_is_being_started
2617    for that violation. It cannot create a deadlock because we are still
2618    running in single threaded mode essentially. Only the IO threads
2619    should be running at this stage. */
```

# double write create

- `buf/buf0dblwr.cc:178, buf_dblwr_create()`

```
173 /****************************************************//**
174 Creates the doublewrite buffer to a new InnoDB installation. The header of the
175 doublewrite buffer is placed on the trx system header page. */
176 UNIV_INTERN
177 void
178 buf_dblwr_create(void)
179 /*==================*/
180 {
181   buf_block_t*  block2;
182   buf_block_t*  new_block;
183   byte* doublewrite;
184   byte* fseg_header;
185   ulint page_no;
186   ulint prev_page_no;
187   ulint i;
188   mtr_t mtr;
189
190   if (buf_dblwr) {
191     /* Already inited */
192
193     return;
194   }
195
196 start_again:
197   mtr_start(&mtr);
198   buf_dblwr_being_created = TRUE;
199
200   doublewrite = buf_dblwr_get(&mtr);
```

# double write create

- `buf/buf0dblwr.cc:178, buf_dblwr_create()`

```
173 /***************************************************************//**
174 Creates the doublewrite buffer to a new InnoDB installation. The header of the
175 doublewrite buffer is placed on the trx system header page. */
176 UNIV_INTERN
177 void
178 buf_dblwr_create(void)
179 /*==================*/
180 {
181   buf_block_t*  block2;
182   buf_block_t*  new_block;
183   byte* doublewrite;
184   byte* fseg_header;
185   ulint page_no;
186   ulint prev_page_no;
187   ulint i;
188   mtr_t mtr;
189
190   if (buf_dblwr) {
191     /* Already inited */
192
193     return;
194   }
195
196 start_again:
197   mtr_start(&mtr);
198   buf_dblwr_being_created = TRUE;
199
200   doublewrite = buf_dblwr_get(&mtr);
```

get double write buffer

# double write create

- `buf/buf0dblwr.cc:178, buf_dblwr_create()`

```
202   if (mach_read_from_4(doublewrite + TRX_SYS_DOUBLEWRITE_MAGIC)
203       == TRX_SYS_DOUBLEWRITE_MAGIC_N) {
204     /* The doublewrite buffer has already been created:
205     just read in some numbers */
206
207     buf_dblwr_init(doublewrite);
208
209     mtr_commit(&mtr);
210     buf_dblwr_being_created = FALSE;
211     return;
212   }
213
214   ib_logf(IB_LOG_LEVEL_INFO,
215     "Doublewrite buffer not found: creating new");
216
217   if (buf_pool_get_curr_size()
218       < ((2 * TRX_SYS_DOUBLEWRITE_BLOCK_SIZE
219     + FSP_EXTENT_SIZE / 2 + 100)
220         * UNIV_PAGE_SIZE)) {
221
222     ib_logf(IB_LOG_LEVEL_ERROR,
223       "Cannot create doublewrite buffer: you must "
224       "increase your buffer pool size. Cannot continue "
225       "operation.");
226
227     exit(EXIT_FAILURE);
228   }
```

doublewrite buffer already created

start of doublewrite buffer creation: skip

# Note

- ## marc_read_from_X
  - read X bytes
  - considering architecture specific endian.

```
175 /***********************************************************//**
176 The following function is used to fetch data from 4 consecutive
177 bytes. The most significant byte is at the lowest address.
178 @return ulint integer */
179 UNIV_INLINE
180 ulint
181 mach_read_from_4(
182 /*=============*/
183   const byte* b)  /*!< in: pointer to four bytes */
184 {
185   ut_ad(b);
186   return( ((ulint)(b[0]) << 24)
187     | ((ulint)(b[1]) << 16)
188     | ((ulint)(b[2]) << 8)
189     | (ulint)(b[3])
190     );
191 }
```

# double write create

- `buf/buf0dblwr.cc:126, buf_dblwr_init()`

```
122 /***********************************************************//**
123 Creates or initializes the doublewrite buffer at a database start. */
124 static
125 void
126 buf_dblwr_init(
127 /*===========*/
128   byte* doublewrite)  /*!< in: pointer to the doublewrite bu
129         header on trx sys page */
130 {
131   ulint buf_size;
132
133   buf_dblwr = static_cast<buf_dblwr_t*>(
134     mem_zalloc(sizeof(buf_dblwr_t)));
135
136   /* There are two blocks of same size in the doublewrite
137   buffer. */
138   buf_size = 2 * TRX_SYS_DOUBLEWRITE_BLOCK_SIZE;
139
140   /* There must be atleast one buffer for single page writes
141   and one buffer for batch writes. */
142   ut_a(srv_doublewrite_batch_size > 0
143       && srv_doublewrite_batch_size < buf_size);
144
145   mutex_create(buf_dblwr_mutex_key,
        &buf_dblwr->mutex, SYNC_DOUBLEWRITE);
```

size of doublewrite buffer
= 2 * EXTENTS (1MB)

# double write create

- `buf/buf0dblwr.cc:126, buf_dblwr_init()`

```
148    buf_dblwr->b_event = os_event_create();
149    buf_dblwr->s_event = os_event_create();
150    buf_dblwr->first_free = 0;
151    buf_dblwr->s_reserved = 0;
152    buf_dblwr->b_reserved = 0;
153
154    buf_dblwr->block1 = mach_read_from_4(
155        doublewrite + TRX_SYS_DOUBLEWRITE_BLOCK1);
156    buf_dblwr->block2 = mach_read_from_4(
157        doublewrite + TRX_SYS_DOUBLEWRITE_BLOCK2);
158
159    buf_dblwr->in_use = static_cast<bool*>(
160        mem_zalloc(buf_size * sizeof(bool)));
161
162    buf_dblwr->write_buf_unaligned = static_cast<byte*>(
163        ut_malloc((1 + buf_size) * UNIV_PAGE_SIZE));
164
165    buf_dblwr->write_buf = static_cast<byte*>(
166        ut_align(buf_dblwr->write_buf_unaligned,
167            UNIV_PAGE_SIZE));
168
169    buf_dblwr->buf_block_arr = static_cast<buf_page_t**>(
170        mem_zalloc(buf_size * sizeof(void*)));
171 }
```

set 2 blocks of double write header

create buffer frame for doublewrite (and make it page sized align)

Prepare buf block array

# Q&A

- Any Questions ?

# Reference

- Source Code : MySQL Community Server 5.6.21

- Transaction Model and Locking
  - [https://dev.mysql.com/doc/refman/5.6/en/innodb-transaction-model.html](https://dev.mysql.com/doc/refman/5.6/en/innodb-transaction-model.html)
  - [http://dev.cs.uni-magdeburg.de/db/mysql/InnoDB-transaction-model.html](http://dev.cs.uni-magdeburg.de/db/mysql/InnoDB-transaction-model.html)