

## 1.background

Nowadays, many cloud databases adopt a storage-computing separation structure, such as Amazon Aurora, Huawei Taurus. The computing nodes run the database program, and the storage nodes load NFS. The network file system exposes the file system interface for computing nodes to use. In this architecture where storage and computing are decoupled, the network becomes a new bottleneck. Buffer Pool. The write amplification caused by frequent write operations in the database has aggravated the severity of the problem. How to optimize network bandwidth and reduce data transmission between nodes is a new challenge facing the database field.

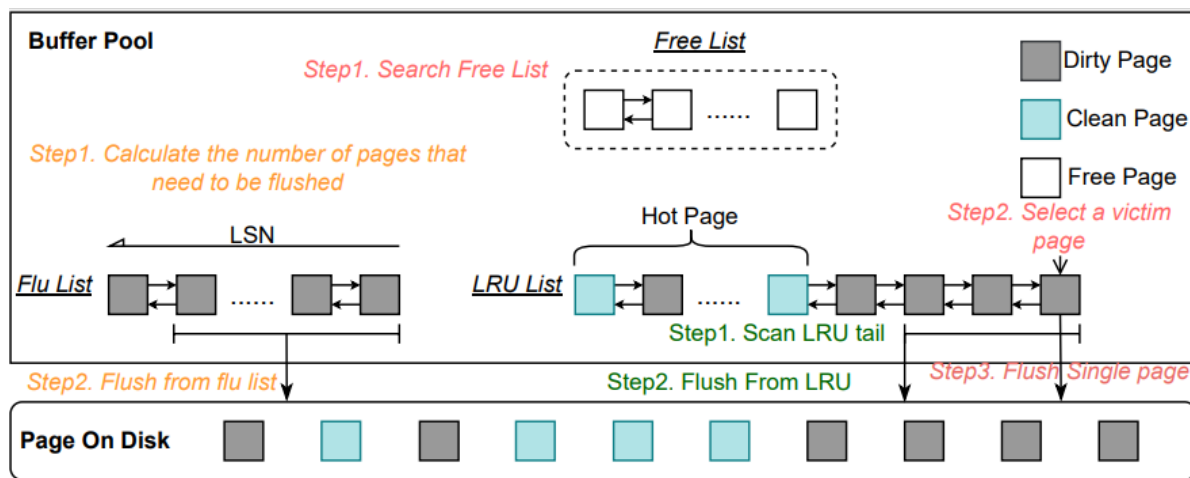
### 1.1 Buffer Pool With write amplification

Buffer Pool. Most OLTP. An important component of the database used to improve database read performance. InnoDB Data is organized into pages and persisted on disk. To alleviate the huge difference in access speed between memory and disk, some data pages are cached in memory and Buffer Pool Manager. These buffer pages are shared and accessed by foreground threads. Buffer Pool Manager. The goal is to cache frequently used pages in memory for as long as possible, while rarely accessed pages are quickly eliminated. In order to maximize the spatial locality of data access, improve the Buffer Pool hit rate, Buffer Pool Manager. The data page is buffered in frames (Buffer Frame) in the form of a linked list and use a cache replacement algorithm (such as LRU, Clock etc.) to manage the buffer frame list. InnoDB Used LRU A variation of the algorithm that keeps the cache hit rate at 90% about.

During database operation, Buffer Pool. Dirty pages are constantly written to disk. In general, Buffer Pool. There are three different types of write operations.

- (1) Since the cache size is limited, in some cases Buffer Pool Manager. Some dirty pages must be written back to disk; (2) In order to prevent dirty page write-back operations from affecting foreground user threads, many databases use a pre-refresh strategy;
- (3) Database system through Checkpoint Mechanism. Cleanup WAL (Write Ahead Log) space, Checkpoint. Will cause a large number of dirty pages to be written.

Next, we will introduce these three different types of page write operations, as shown in the figure below. 1. The process of three types of page write operations is shown.



picture1. Buffer Pool Write Operation

The data page in memory is called a buffer frame. The buffer frame not only maintains a copy of the disk page, but also maintains the metadata information required by the disk page in memory, such as reference counts. The buffer frame has three states in memory:

- (1) Free: The buffer frame is free and has no stored disk pages.
- (2) Clean: The buffer frame stores a copy of the disk page, but the page has not been modified by any transaction.
- (3) Dirty: The buffer frame stores a copy of the disk page that has been modified by at least one transaction.

#### 1.1.1 BUF\_SINGLE\_PAGE

Buffer Pool Manager. Organize the free buffer frames into a linked list (Free List), when the page requested by the current transaction encounters Cache Miss. When Free List. Are there any idle buffer frames in **Step 1**. Otherwise, when Buffer Pool. When there are no free buffer frames available, it is necessary to scan from back to front according to the cache replacement algorithm LRU List. Find the status Clean Buffered frames (**Step 2**). Once a state is encountered Clean. If the buffer frame is set, the disk page can be read into the buffer frame. Otherwise, LRU. The tail dirty pages are flushed to disk synchronously.

(**Step 3**) to get a free frame to use, and then read the missing page into the free frame.

This operation is called `BUF_SINGLE_PAGE`, when Buffer Pool Full and no status is Clean Idle frame, it must comply with Read-After-Write (RAW). This is an expensive operation, and the foreground transaction must pause and wait for the dirty page synchronization write to complete before issuing a read request. Buffer Pool The utilization rate is very low. When the foreground transaction is busy, each page read will trigger RAW. This will cause serious performance degradation and is unacceptable for the database system.

picture1 The red text shows `BUF_SINGLE_PAGE` process.

### 1.1.2 BUF\_LRU\_PAGE

To alleviate RAW, the impact caused by this is that modern databases usually use Preflush/Pre-refresh mechanism, the database system background will have Dirty Page Writer Thread, this page refresh method is called `BUF_LRU_PAGE`. InnoDB uses Redo Log (WAL) To ensure the persistence of committed transactions, use Undo Log To ensure the atomicity of uncommitted transactions, the refresh of data pages can be done by Steal/No-Force strategy, data pages can be written to disk at any time during transaction execution without affecting the database ACID Features. The background thread periodically LRU Tail of linked list (**Step 1**), using asynchronous IO Interface such as Linux AIO Write a certain number of dirty pages to disk (**Step 2**), thereby reducing the triggering of foreground transactions RAW. The number of operations.

picture1 The green text shows `BUF_LRU_PAGE` process.

### 1.1.3 BUF\_FLUSH\_PAGE

Database systems use write-ahead logs (WAL) Stores any changes made by the transaction to the database to ensure the durability of the transaction. When the database crashes, only the WAL, apply it to the database system, and you can restore the state before the crash. Usually the recovery time is proportional to the log size. In order to keep the recovery time short, you need to execute it regularly. Checkpoint Forcefully flush dirty pages to truncate the log. This page flushing method is called `BUF_FLUSH_PAGE`.

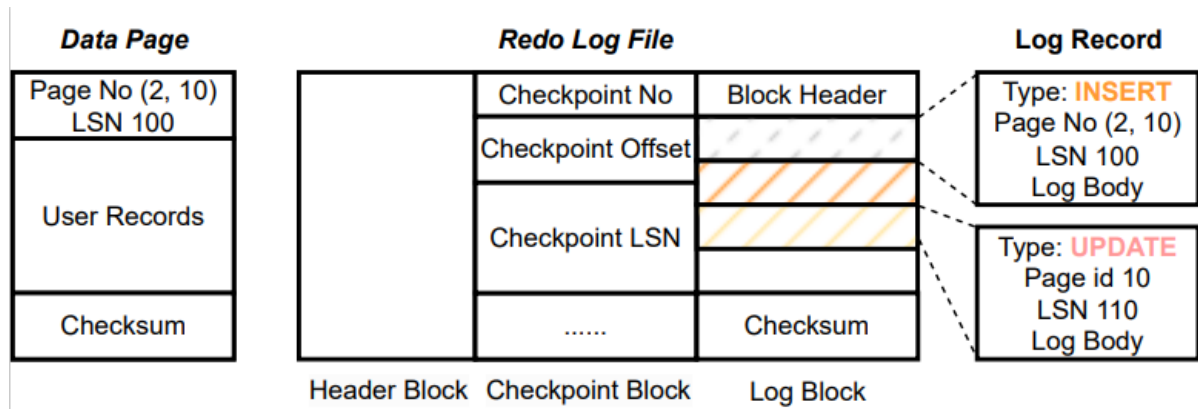
InnoDB The dirty buffer frame is LSN Organized from large to small FLU List, when the background thread executes Checkpoint When the log is generated, the number of logs to be generated will be calculated based on the current log generation speed and the remaining log capacity. FLU List Number of dirty buffer frames flushed  $N$  (**Step 1**), in general, the faster the log is generated,  $N$  The larger the value of . Then the dirty buffer frame is flushed to disk (**Step 2**).

picture1 The orange text shows `BUF_FLUSH_PAGE` process.

## 1.2 Redo Log With crash recovery

because Buffer Pool With the existence of transactions, the modification of disk data by transactions may lag behind the memory. If the database process or machine crashes, the memory data will be lost. Most modern databases use ARIES The algorithm or its variants ensure the durability and atomicity of transactions, and the database strictly adheres to Write Ahead Logging (WAL) rule, all changes made by a transaction to a data page need to be recorded in the write-ahead log before being reflected on disk, and the log must be written before the corresponding Page. When a failure occurs and causes memory data to be lost, the database can be restarted by replaying the pre-written log. Page Restore to the state before the crash. InnoDB The write-ahead log in Redo Log.

As shown below 2 As shown, InnoDB Taken Physiological Logging Ways to organize Redo Log, by Page Record logs for units, and the logs record Page Logical changes made. All logs are Log Sequence Number (LSN) to uniquely identify the database. Page The header information will record the last modification log of the page. LSN. Depending on the modification operations performed by the firm, Log Record There are different types of Page of User Record When a record is inserted into INSERT Type of log, while updating Page The corresponding records in UPDATE Type of log. Log Body The part stores the data needed for log replay. The content stored in this part varies according to the log type. In order to achieve the expected crash recovery time, the database system background will continuously perform Checkpoint, 1) from FLU List Refresh a certain number of pages, we assume that the largest of these pages LSN for Checkpoint LSN; 2) Checkpoint LSN Will be recorded Log In the document, this indicates that Checkpoint LSN All previous logs can be truncated, and their corresponding modifications have been persisted to disk.



picture2. Redo LogStructural diagram

When the database is restarted, it scans Redo Log to perform crash recovery, which is divided into two stages.

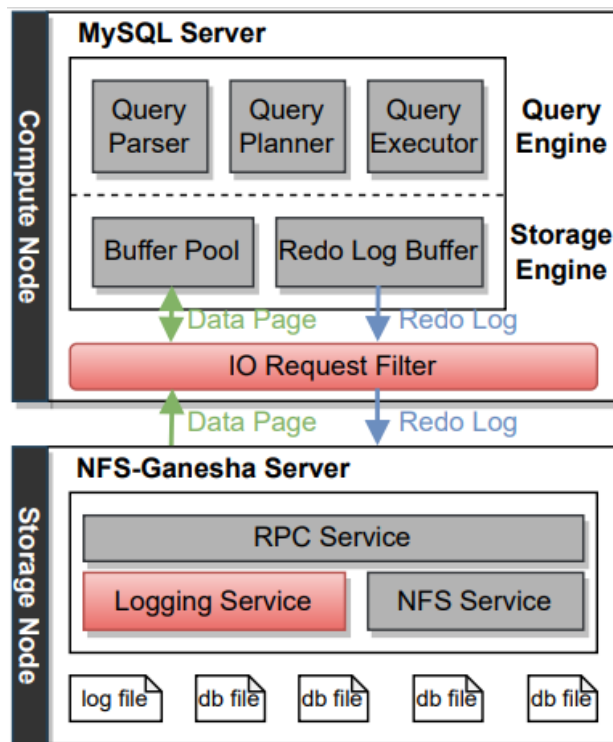
## 1.2.1 Redo stage

This phase is responsible for restoring the database to the state before the crash. This phase does not distinguish between the committed states of transactions. Modifications made by both committed and uncommitted transactions will be restored. Checkpoint LSN: The changes corresponding to the previous logs have been persisted, so this stage will only check the LSN of the current log record and the LSN of the page. The log is reapplied to the page only when the size is smaller.

## 1.2.2 Undo stage

InnoDB: The changes made by uncommitted transactions are stored in Undo of Page middle. Undo Page: Also affected Redo Log of protection, so in Redo stage will use Undo Page to roll back transactions that were still active at the time of the crash.

# 2. Solution



picture3. log pushdown Overall structure

MySQLAs a commonly used relational database management system, it plays an important role in cloud database systems. However, while the storage-computing separation architecture improves scalability and availability, it introduces the problem of network bottlenecks. In addition,Buffer PoolandWALThe mechanism leads to serious write amplification, which has a significant impact on the performance of the entire system. To address these challenges, we propose a non-intrusive improvement solution to reduce networkIOAnd reduce write amplification.

TraditionalMySQLThe log system requires that all data modification operations be persisted toRedo Logand flush it to disk. We noticed that we only needed to ensureRedo LogBased on this observation, we propose a non-intrusive solution that bypassesBuffer Poolmiddle Data PageThe down-brush operation can reduce the networkIO, and use the storage nodeCPUResources implement log replay to reduce the impact of write amplification.

Our non-invasive approach offers the following advantages:

- AvoidMySQLSource code modification: Our solution does not require modificationMySQLThis means that no additional risk and complexity are introduced. It can be used with existingMySQLDeployment is compatible and easy to implement and maintain.
- Reduce networkIO: By swiping down onlyRedo Logand bypassData PageBy refreshing the network, we can greatly reduce theIOThis is particularly important in the storage-computing separation architecture, which can significantly reduce the consumption of network transmission bandwidth and improve the performance and response speed of the entire system.
- Take full advantage of storage nodesCPUResources: Our solution relies on the storage nodesCPUThe system uses the resources to replay logs instead of relying on the speed of network transmission. This allows log replay to be completed quickly, avoiding delays in reading and waiting, and further improving system performance and availability.

Through this non-invasive improvement solution, we can solve the problem ofMySQLThe network bottleneck and write amplification problem of the log system. Our approach brings better performance and scalability to the cloud database system, providing users with more efficient and stable data management services.

We are based onMySQL 5.7.30andNFS-GaneshaA prototype system was implemented, calledLogPushDown,picture3Shown LogPushDownthe overall structure.

NFS-GaneshaIt is a widely used program that runs in user mode.NFSService, which implements the completeNFSProtocol, which can provide a file system interface to the outside world.NFSService, it is more flexible and more conducive to the development of prototype systems.

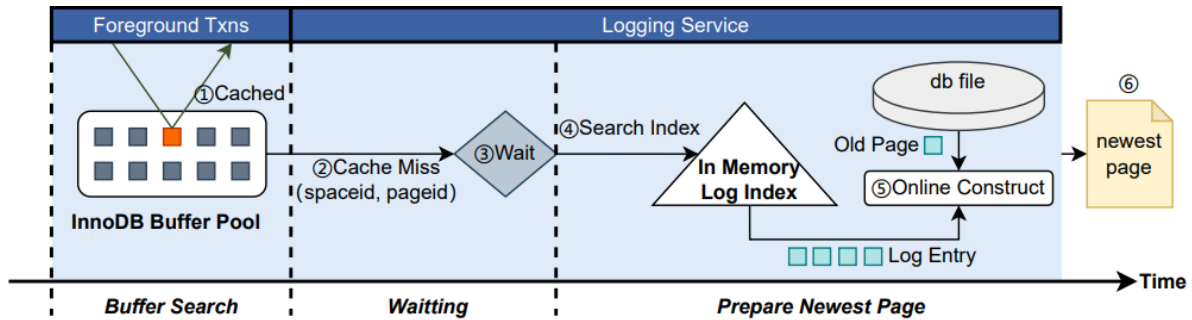
Although the prototype system we implemented is specific toMySQLdatabase, but consideringRedo LogThis is the solution adopted by most disk-oriented database systems, so we believe that this method is universal and can be transplanted to other database systems.

**IO Request Filter.**To bypassBuffer PoolThe refresh operation of the data page is incorrect.MySQLTo prevent any impact on applications, we implemented a pluggableIORequest filter,MySQLThis component is unaware of theMySQLIssuedData PageWrite requests are filtered out, keeping onlyRedo LogWrite request.

**Logging Service.**whenBuffer PoolWhen the cache misses, the correspondingData Page,Logging Service WillRedo Log Manage and replay logs regularly to ensureData Page ReadRequest to read the latest version of the page.

## 2.1Data page read path

Below4It shows the path of the foreground transaction when requesting the data page. It can be divided intoBuffer Search,Waiting, Prepare Newest PageThree stages.



**Buffer Search.** If Buffer Pool Target in cachePage, it will be returned to the foreground transaction immediately ①, otherwise it will generate Cache Miss, Buffer Pool The corresponding page will be requested from the storage server ②).

**Waiting.** Data Page The read request will be stored in the node Logging Service deal with. Logging Services A memory log index is maintained, and there will be Log Parser The thread continuously parses the log file and updates the index. The parsing of the log file may lag behind the speed of log file generation. Data Page When a read request comes, Log Parser The latest log files may not be parsed and updated to the index structure in time, and the unparsed logs may contain Page Related logs. We need to record the latest log generated when the request arrives LSN, then spin wait log parser Analyze the latest log ③).

**Prepare Newest Page.** Logging Service There will be corresponding Apply Worker The thread periodically performs log replay operations. Apply Worker The thread will select some logs from the log index and replay them Data Page The replayed logs will be deleted from the log index. Data Page When a read request comes in, the log index is retrieved. Page If the relevant log ④ cannot be retrieved, it proves that Page Has been Apply Worker The thread is updated to the latest version. At this time, no other operations are required, just return the latest page ⑥). Otherwise, we need to replay the log online to build the latest version. Page ⑤)

### 3.Implementation details

## 3.1Non-invasive IO Filters

MySQL All IO Requests are made using POSIX Standard file system call interface, such as open, pread, pwrite Etc. These system calls are usually libc In the realization, LD\_PRELOAD The mechanism allows the specified shared library to take precedence over the system default library. Therefore, when a function in a program is called, the dynamic linker will first look for the corresponding function implementation in the preloaded shared library, and if it is found, it will use the preloaded function implementation. Therefore, we will IO The filter is implemented as a shared library libcatcher, libcatcher The same name is implemented in pread, pwrite Wait for system call, replace libc Default system call behavior.

Specifically, libcatcher The following system calls are intercepted: open, open64, pread, pread64, pwrite, pwrite64, close.

Of course, in some cases we still want to use libc system call, dlsym The function can help accomplish this task. libcatcher The library is loaded using dlsym, libc Find the address of the system call with the specified name in the library, register it, and use it later.

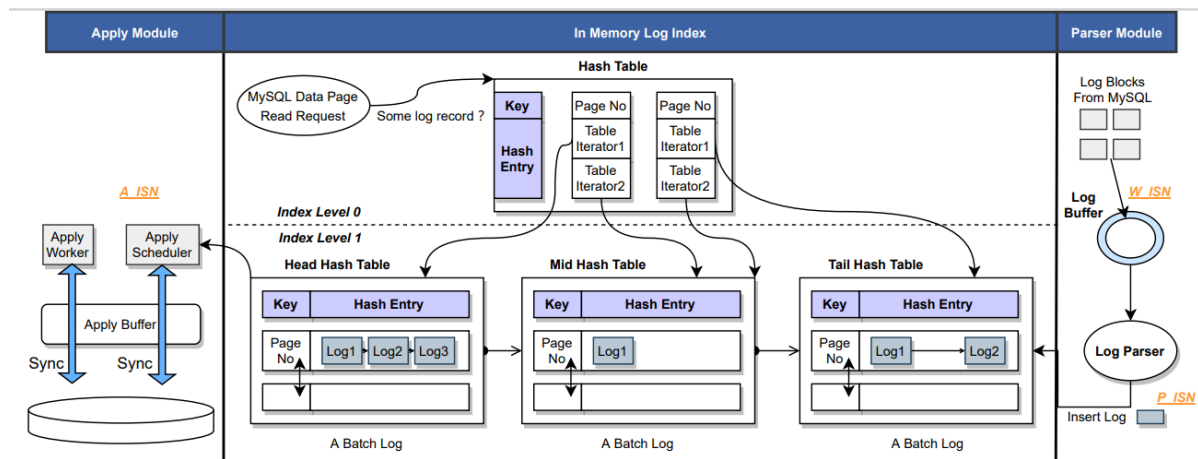
MySQL middle Redo Log The file has a specific file name, and the format of the database file is \*.ibd, therefore, in open In the system call, different files can be distinguished according to the file name. If it is a database file (ibd), we add its file descriptor to the filter list, so that, in pwrite The file descriptor filter table can be searched in the system call to directly filter the write operation of the database file.

This method does not require MySQL To make any changes, just MySQL At startup, via LD\_PRELOAD System variable assignment libcatcher Location, such as:

```
LD_PRELOAD=./libcatcher.so ./mysqld
```

This allows for seamless integration of filter components.

## 3.2 Logging Service



picture5. Logging service

### 3.2.1 Parser Module

When MySQL writes logs to a storage node, they are not only written to the log file on disk, but also to the log file in memory. Log Buffer is a copy of the log file in memory. Log Buffer is a circular buffer that only caches the latest unplayed logs. Its size is Batch Size. Through experiments, we found that in the scenario of storage and computing separation, due to the influence of network bandwidth, lock contention, etc., MySQL transaction system handles transaction generation Redo Log. The speed is always slower than the speed of log playback. In our design, since the background thread will batch (Batch Size) Playback Redo Log, so in the worst case, there will be a batch of unfilled Batch Size log is receiving write requests, and a batch is full Batch Size. The log is being replayed, so we will Ring Buffer. The size is set to 2 times Batch Size is reasonable. In the future, if computing nodes use more powerful CPU, while using higher bandwidth network technologies, such as RDMA, there may be Redo Log. The generation speed is faster than the log playback speed. As the system runs, the logs that are not played back are Redo Log. Will gradually accumulate, on the one hand Log Buffer Dynamic expansion may be required. Redo Log will accumulate infinitely, requiring a large amount of memory space. On the other hand, the unplayed data accumulated in the storage node Redo Log. The more, MySQL Reading data pages requires Online Construct. Therefore, users need to give Log Buffer Set a reasonable upper limit when Redo Log. The accumulation reaches the threshold and must be paused MySQL Write Redo Log.

It should be noted that in order to avoid Redo Log Part of the write (Partial Write) question, MySQL by Block (512Byte) Redo Log, and Log Entry. The size of is arbitrary, as shown in the figure 2 middle Redo Log. The file structure is shown below, except Log Record. In addition, Log Block Some additional metadata is stored in the header and footer. Ring Buffer Only the Log Record, and does not store this metadata. Log Record Not enough to fill the entire Log Block, the rest will be 0Byte padding. MySQL Two consecutive times Redo Log. In a write request, there may be two Log Block. The situation of repeated writing, so we need to solve two problems, 1) right Log Block Remove the metadata at the beginning and end. 2) and accurately extract Log Record The incremental part of .

We use three pointers to manage Log Buffer Space: W\_ISN, P\_ISN, A\_ISN. ISN (internal sequence number) is an integer that increases globally monotonically according to the number of bytes.

**W\_ISN** We maintain a W\_ISN to indicate where the next write request should be written. NBytes in size Redo Log Record. When writing, W\_ISN will grow accordingly Nbytes.

**P\_ISN** From the figure 5 Be able to see. Log Parser Will continue to parse the log, each parsing NBytes of log, P\_ISN Will grow accordingly N bytes. P\_ISN = W\_ISN hour Log Parser Will pause parsing and wait Log Buffer Write a new log in.

**A\_ISN** From the figure 5 Be able to see. Apply Worker The logs will be played back continuously. Once a batch of logs are played back, A\_ISN Will move forward Batch Size To prevent unplayed logs from being overwritten by newly written logs, only when W\_ISN - A\_ISN < Log Buffer Size Log writing is allowed only when

MySQL There are two Redo Log Files, which form a log file group, are used in a circular manner. Log Buffer It is much smaller than the log file group, so the third problem needs to be solved. 3) How to Log Blocks Extracted from Log Record Map to a smaller range Log Buffer The code snippet shows how to Redo Log blocks Write to Log Buffer process.

### 3.2.2 In Memory Log Index

As shown5As shown, we use the hash table asRedo LogMaintains a two-level memory index to speed upRedo LogWe will first1The layer index is organized into a hash table array. The hash table isPage Nois the key and the value is aLog Record's linked list.

Each hash table in the first-level index has three states:Open,Close,Applying.

The hash table at the end of the queue isOpenstate,Log ParserSingle-threaded sequential parsingLog BufferAccording to the log inLog RecordIn Page Nofield, will belong to aPageThe log is added to the hash table at the end of the queue.Log ParserThere are two main reasons for designing it as a single thread instead of using multithreading for acceleration: first,MySQLofRedo LogThere is no field in theLog Recordlength, so for aLog RecordFor example, all fields must be parsed sequentially to know the nextLog RecordThe location to which it belongs; secondly, to ensure the correctness of log playback,LSNSmall logs need to be replayed first, so theLog RecordThe linked list must beLSNArranged in order, multiple threads parse outLog RecordNo guaranteeLSNMonotonically increasing order. For these two reasons,Log ParserMust be designed to be single-threaded. OpenThe total number of log bytes stored in the hash table of the state reaches a certain threshold (Batch Size) will be converted toCloseStatus, no longer receivingLog ParserWrite request.Log ParserAn empty hash table will be reinitialized and added to the end of the hash queue to receive new write requests. The hash table at the head of the queue isApplyingstatus, where the logs will beApply ModuleConsume.

exist3.2.1In the previous section, we mentioned that when the log generation rate is greater than the log playback rate, logs will accumulate, which will inevitably cause the hash queue of the first-level index to be too long.MySQLSend a page read request,Logging ServiceNeed to be doneOnline ConstructWhen the first hash table is scanned, all hash tables in the hash queue need to be scanned. In the worst case, there are no related logs in all hash tables, which means this is an invalid scan. To speed up this process, we build an additional hash table as the first0Layer index, the hash table isPage NoThe key is the value, which is a hash table pointer array pointing to the1The layer index stores thePage NoThe hash table of the relevant logs. The data page read request only needs to query the hash table to find the corresponding log.1The position in the layer index greatly improves the query efficiency.

### 3.2.3 Apply Module

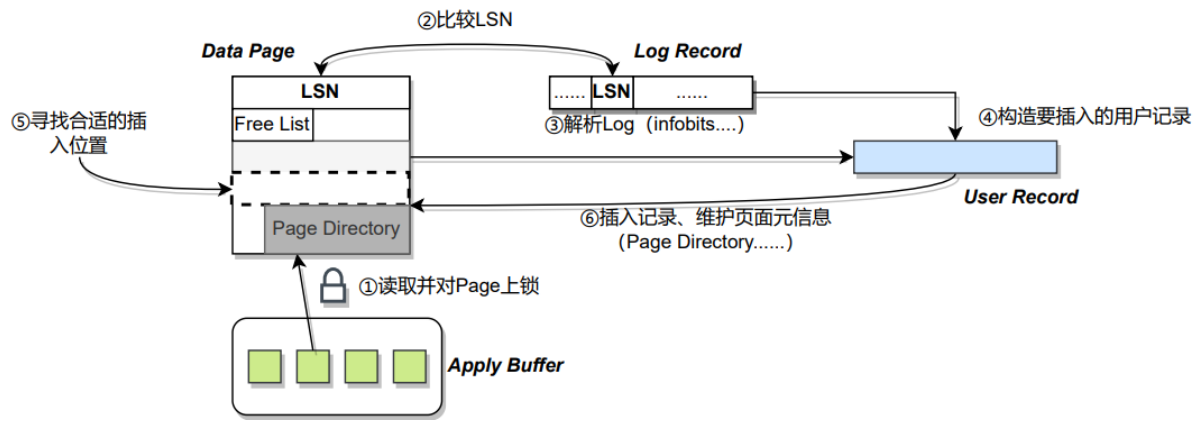
Apply ModuleProvides a "best effort" service that continuously replays logs from the index in the background in an effort to build the latest version of the page. MySQLThere are as many as61There are three types of logs. According to the type of logs acting on the target page, they can be divided into6Types, as shown in the following table1As shown. Among them, the effect isRTree pages and logs acting on compressed pages are optional.RTree the tree isMySQLThe data type used to represent spatial data in . This type of log is generated only when spatial index is enabled. Similarly, logs for compressed pages are generated only when page compression is enabled.MySQLUse a separateUndoType of page to storeUndoLogs also need to record the corresponding types of logs to ensureUndo Page persistence.MySQLCorrespondingCompactandRedundantThere are two row formats to format the row records in the user table. Different row formats require different information to be recorded during crash recovery and playback, so different types of logs are required.MySQL 5.0Start by usingCompactWhen creating a file or modifying the file name, the file system meta information (inode), in order to prevent the loss of metadata due to power failure, the corresponding logs must also be recorded.

作用于文件信息	作用于 B+ 树页面		作用于 R 树页面	作用于 Undo 页面	作用于压缩页面	其它
FILE_NAME、	Compact 行格式	Redundant 行格式	PAGE_CREATE_RTREE、	UNDO_INSERT、	ZIP_PAGE_COMPRESS、	TRUNCATE、
FILE_RENAME、	COMP_REC_INSERT、	REC_INSERT、	PAGE_COMP_CREATE_RTREE	UNDO_HDR_REUSE、	ZIP_PAGE_WRITE_HDR、	INDEX_LOAD、
FILE_CREATE	COMP_PAGE_REORGANIZE	LIST_END_DELETE		UNDO_HDR_CREATE	ZIP_PAGE_REORGANIZE	WRITE_STRING

surface1. Redo Logtype

Each type of log has different playback logic, and it is very laborious to fully support the playback of all log types.TPC-C, SysbenchThe workload generatedRedo LogWe found that it will only affect the file information andB+Tree Page (Compactrow format), which acts on UndoPages and some other logs, in which theB+TreeThe proportion of logs of tree page type has reached95%, so we only implementB+Tree page type logs and some other necessary types of logs, a total of17ForUndoPages and other unsupported log types, we choose to comply withMySQLThe original logic,IO Filter The swipe down operation of the corresponding page is allowed in the component.

As shown below6WeB+Insert a record into the treeCOMP\_REC\_INSERTThis article takes the log of the type as an example to explain how to perform log playback.

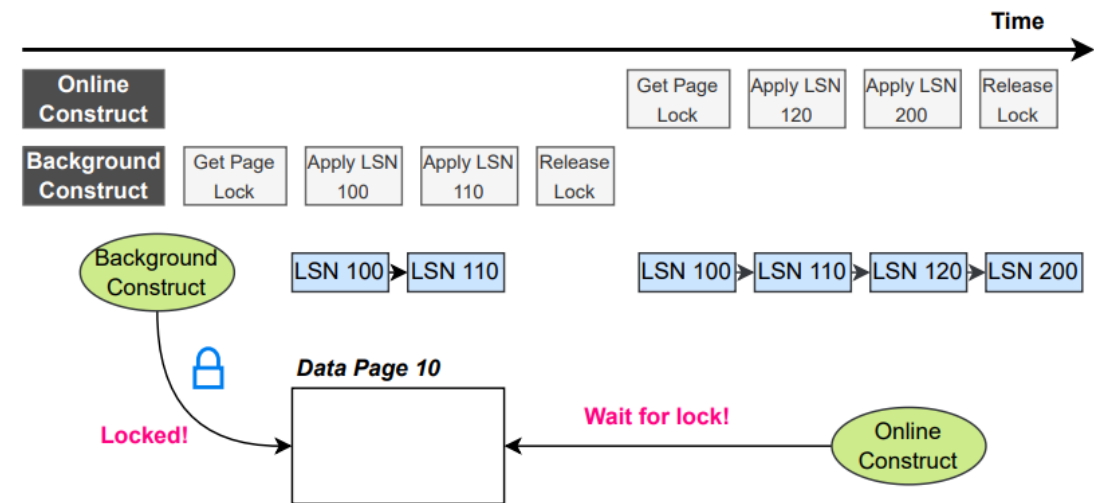


picture6. COMP\_REC\_INSERTLog playback process

from Apply Buffer Read the page and lock it (①). Log playback requires that the page must be in the correct state, so a log playback operation must be atomic. Any thread that wants to play back a page must first obtain the page lock before proceeding. Redo Log of LSN and Data Page of LSN (②), if Redo Log LSN Less than Data Page LSN, which proves that the data page has been replayed to a newer version by other threads. This log replay can be skipped directly.

As shown below 7, imagine a scenario where the background playback thread starts Level 1 In Memory Log Index of Head Hash Table Extracted from LSN for 100 and 110. Two logs, the foreground process from the entire In Memory Log Index Extracted from LSN for 100, 110, 120, 200 of four logs, all of which they want to Page 10 To replay the log, they need to compete for the page lock. Assuming that the background thread first competes for the lock, it replays the log. 100, 110. After that, the lock is released, and the foreground thread can obtain the lock of the page. LSN Discovery, Log 100, 110. It will skip these two logs and only replay the log 120 and 200.

In order to minimize Log Record The size of COMP\_REC\_INSERT This type of log will only record the inconsistent parts with the previous user record, which needs to be parsed Log Record (③), get the previous user record in Data Page Position in, merge Log and the data of the previous user record to construct a complete user record (④), and then you need to Page Apply for a space in (⑤) (may be reused) Free List space), copy the newly constructed user record to the specified space (⑥), and at the same time, maintain the meta information in the page, the most important of which is to update the page LSN.



picture7. Online Construct and Background Construct Conflict

Only by speeding up the background log playback rate, MySQL of Buffer Pool produce Cache Miss Before reading the data page, replay the logs related to the data page as early as possible to reduce Online Construct. The number of log entries that need to be replayed during the process can avoid page read requests pausing for too long, thus affecting transaction throughput and latency. Apply Module Designed as a Scheduler and multiple Worker. The multi-threaded model. Head Hash Table Transformed into Close. When the state is reached, a round of background log replay will be started immediately. Scheduler Distribute the logs in the hash table evenly to several buckets. Worker, logs in the same hash bucket will only be assigned to one Worker, to reduce contention between threads. Worker After completing the log playback of a page, the page needs to be written back to disk immediately. Scheduler After the tasks are assigned, it will also become Worker Perform log playback and wait for all worker After completing the work and becoming idle, Scheduler Will move forward A\_ISN.



## 4.Environment Construction

### 4.1 MySQLend

This article is based onMySQL 5.7.30To build, you first need to buildMySQLenvironment.

1.First install the dependencies:

```
sudoapt install pkg-config libssl-dev bison-y
```

2.Configuration GenerationMakefile:

```
# Enter the project root directory
mkdircmake-build-release cmake
-S. \
-Bcmake-build-release \
-G"Unix Makefiles"\
-DCMAKE_BUILD_TYPE=Release \
-DCMAKE_INSTALL_PREFIX=path/to/mysql \
-DMYSQL_DATADIR=path/to/mysql/data \
-DMYSQL_UNIX_ADDR=path/to/mysql/data/mysql.sock \
-DSYSCONFDIR=path/to/mysql/data \
-DWITH_DEBUG=ON \
-DWITH_BOOST=path/to/mysql-code/boost \
-DMYSQL_MAINTAINER_MODE=OFF
```

3.Compile and generate executable files:

```
# - - - j 16Specifies the number of threads to use when compiling the project.
cmake--buildcmake-build-release--targetall-- -j16
```

4.At startupMySQLBefore that, you need to initialize the data directory and setrootAccount permissions and password:

```
# Create the data directory before initializationMySQLUser groups and
MySQLuser sudo groupadd mysql
sudo useradd -r -gmysql -s/bin/false mysql
```

5.initializationmysqldData directory, enter the project root directory, and perform the following operations:

```
cd./cmake-build-release/sql
.mysqld--basedir=path/to/mysql--datadir=path/to/mysql/data--lower_case_table_names=0--user=
mysql--innodb-flush-method=O_DIRECT-- innodb_flush_log_at_trx_commit=1--innodb_log_file_size=
2G-- innodb_change_buffering=none--default-storage-engine=InnoDB--default-tmpstorage-engine=
InnoDB--disabled_storage_engines=MyISAM--innodb-checksumalgorithm=none--
innodb_log_checksums=OFF--innodb_doublewrite=0--initializeinsecure
```

6.Solve startup **Error** :

MySQL Dependency when printing error logs errmsg.sys File, compile MySQL Afterwards, will be cmake-build  
debug/sql/share Generate various language versions in the folder errmsg.sys file, copy the file to  
/home/lemon/mysql/share Folder.

```
mkdir path/to/mysql/share/  
cp path/to/mysql-code/cmake-build-release/sql/share/english/errmsg.sys mysql/share/ path/to/
```

7. Add read, write and execute permissions to the data directory:

```
chmod-R777 path/to/mysql/data
```

8. start up Server end mysql:

```
.mysql --basedir=path/to/mysql --datadir=path/to/mysql/data --  
socket=path/to/mysql/data/mysql.sock --lower_case_table_names=0 --user=mysql --innodb-flush-method  
=O_DIRECT --innodb_flush_log_at_trx_commit=1 --innodb_log_file_size=1G --innodb_change_buffering=  
none --default-storage-engine=InnoDB --default-tmp-storage-engine=InnoDB --disabled_storage_engines=  
MyISAM --innodb-checksum-algorithm=none --innodb_log_checksums=OFF --innodb_doublewrite=0
```

9. start up Client end mysql: Reviser root User password, enter the project root and record it, and perform the following operations:

```
/path/to/mysql-code/cmake-build-release/client/mysql -uroot -h127.0.0.1 -P3306 -p
```

10. No need to enter a password, just press Enter to start MySQL Client login MySQL. After that, modify the code on the client side and close it. Server Duanhe Client end:

```
ALTER USER 'root'@'localhost' IDENTIFIED WITH mysql_native_password BY "root"; FLUSH privileges;  
  
shutdown;  
exit;
```

## 4.2 Preparing test data

bysysbench This article takes the benchmark test tool as an example to introduce how to generate test data.

Restart MySQL, use the client tool to connect and create a database for testing:

```
CREATE DATABASE sbtest;
```

Then use sysbench to initialize test data:

```
sysbench --threads=20 \
--mysql-host=127.0.0.1 \
--mysql-port=3306 \
--mysql-user=root \
--mysql-password=root /usr/share/sysbench/
oltp_common.lua \
--tables=40 \
--table_size=200000 \
prepare
```

## 4.3 NFS-Serverend

1.First install the dependencies

```
sudo apt install g++ libboost-dev cmake git doxygen sudo apt install build-essential
libglu1-mesa-dev libc6-dev sudo apt install libkrb5-dev libsasl2-modules-gssapi-mit
libkrb5-dev sudo apt install liburcu-dev
```

```
sudo apt install libcap-dev libtirpc-dev rpcbind sudo apt-get install
uuid-dev libacl1-dev liblzo2-dev
```

Note: Startup `nfs-ganesha` Before you install `nfs-kernel-server` Otherwise, when you start `nfs-ganesha` When you get  
When I get the following error, I guess some dependency might be missing? I found that it should be missing `rpcbind` .

Cannot register NFS V3 on TCP

2.Install `nfs-kernel-server`

```
sudo apt install nfs-kernel-server
sudo /etc/init.d/nfs-kernel-server stop #close nfs-kernel-server Because we don't need to use it
```

3.Revise NFS-Server Configuration Files:

turn up `nfs-ganesha-sql/src/include/applier/config.h` File, modify `LOG_PATH_PREFIX`, `DATA_FILE_PREFIX`, `DATA_FILES`,  
`PER_FILE_LOG_FILE_SIZE` Four variables, among which `DATA_FILES` Variables are what you create `sysbench` All table data files,  
`LOG_PATH_PREFIX` yes Redo Log The file path where it is located, `DATA_FILE_PREFIX` yes MySQL The main directory of the  
database table data files, `PER_FILE_LOG_FILE_SIZE` Every Redo Log The size of the log file.

4.Compile the source code:

```
# Enter the project root directory
cd nfs-ganesha-sql
mkdir src/build
cmake -S src -B src/build -DCMAKE_BUILD_TYPE=Release -DUSE_FSA_VFS=ON -DUSE_9P=OFF
-DUSE_FSA_LUSTRE=OFF -DUSE_FSA_GPF=OFF -DUSE_NLM=OFF -
DUSE_FSA_CEPH:STRING=OFF -DUSE_FSA_GLUSTER:STRING=OFF -
DUSE_FSA_KVSFS:STRING=OFF -DUSE_FSA_MILLIARDF:STRING=OFF -
DUSE_FSA_PROXY_V3:STRING=OFF -DUSE_FSA_MEM:STRING=OFF -
DUSE_FSA_RGW:STRING=OFF -DUSE_FSA_S3:STRING=OFF -
DUSE_FSA_S3RBW:STRING=OFF -DUSE_FSA_S3RBW_V4:STRING=OFF -
cmake --build src/build --parallel 20 -DUSE_GSS=OFF
```

## 5.start upnfs-server:

Create a configuration file first

```
vim~/ganesha.conf
```

The configuration file content is as follows, you need to modify `Plugins_Dir` parameter, so that it points to `nfs-server` Compile the generated library directory, `Path` The parameter represents the file directory you need to expose to the client.

```
NFS_CORE_PARAM {
    Plugins_Dir =    path/to/nfs-ganesha-sql/src/build/lib;
}
EXPORT
{
    Export_Id =    77;
    Path =    path/to/mysql-Dir;
    Pseudo =    /;

    FSAL {
        Name = VFS;
    }
    Access_Type    = RW;
    Disable_ACL    = true;
    Squash =    No_Root_Squash;
    Protocols    = 3,4;
}
EXPORT_DEFAULTS
    Transports = UDP, TCP;
    SecType = sys;
}
```

`no_root_squash`: Indicates that when the client `root` Local access when identity is granted `root` Permissions (default is `root_squash`).

`root_squash`: Indicates that the client uses `root` When a user accesses the shared directory, `root` The user is mapped to an anonymous user. Create a

directory to store `ganesha` of pid information:

```
sudo mkdir/var/run/ganesha
```

The final privilege escalation start `nfs` server:

```
sudo path/to/nfs-ganesha-sql/build/bin/ganesha.nfsd -F -L /dev/stdout -f path/to/ganesha.conf -N NIV_EVENT
```

- `F` (foreground), Will `ganesha.nfsd` The service runs in the foreground, - `L` Specify the location where the log will be output. `f` Specify the path where the configuration file is located. `N` Specify the output log level.

## 4.4 Start and test

1. Client `MySQL` One side needs to be mounted `NFS` service, the server `mysql` Directory exposed to `mysql` One side:

```
sudo mount -t nfs4 11.11.11.12:/ path/to/mysql-Dir (Local mount point)
```

2. Revise `io_filterComponent`, compile, open `libcatcherProject`, again `catcher_filter.cpp` Modifications in the file `data_file_setVariable`, pointing to `allsysbench` Generate the database table file, and then compile it, the compiled `.so` Dynamic library is called `libcatcher_filter.so`.
3. start up `MySQL`, in `LD_PRELOAD` specify `libcatcher_filter.so` The location of the dynamic library, `--basedir` and `datadir` should be the remote directory you mounted.

```
LD_PRELOAD=path/to/libcatcher/libcatcher_filter.so ./mysqld --basedir=path/to/ mysql --datadir=path/to/mysql/data --
socket=/tmp/mysql.sock --lower_case_table_names=0 --user=mysql --innodb-flushmethod=O_DIRECT --
innodb_flush_log_at_trx_commit=1 -- innodb_change_buffering=none
-- default-storage-engine=InnoDB -- default-tmp-storage-engine=InnoDB --
disabled_storage_engines=MyISAM -- innodb-checksum-algorithm=none --
innodb_log_file_size=2G --innodb_buffer_pool_size=128M --innodb_doublewrite=0 --
innodb_log_checksums=OFF
```

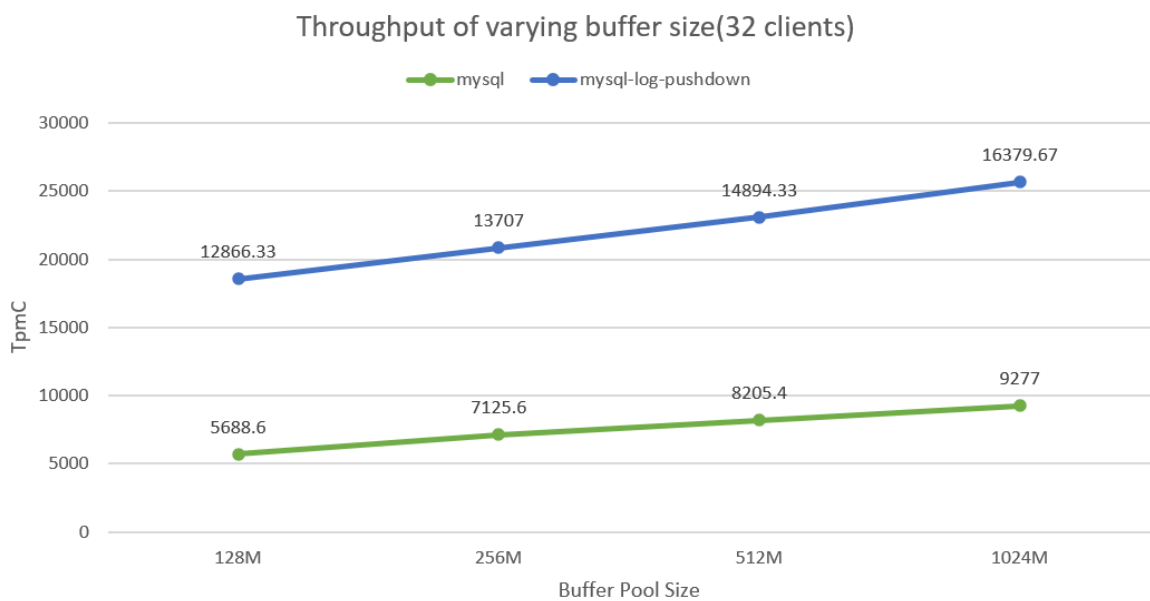
#### 4. start `sysbench` To test:

```
sysbench --threads=8 \
--time=180 \
--report-interval=10 \
--mysql-host=127.0.0.1 \
--mysql-port=3306 \
--mysql-user=root \
--mysql-password=root /usr/share/sysbench/
oltp_read_write.lua \
--tables=40 \
--table_size=2000000 \
--mysql-ignore-errors=2013 \
--mysql-socket=/tmp/mysql.sock \ run
```

## 5. Experimental results display

We built a storage and computing separation system to conduct experiments and test `LogPushDown` The storage node has the same configuration as the computing node, with an additional `PCIe3.0 Samsung 980` Solid-state drives are used to store database files. `RDMA` Network card (still using the normal `TCP/IP` Protocol) for network communication, and computing nodes run `MySQL-Server` and `TPC-C` Workload, storage node running `NFS-Ganesha`.

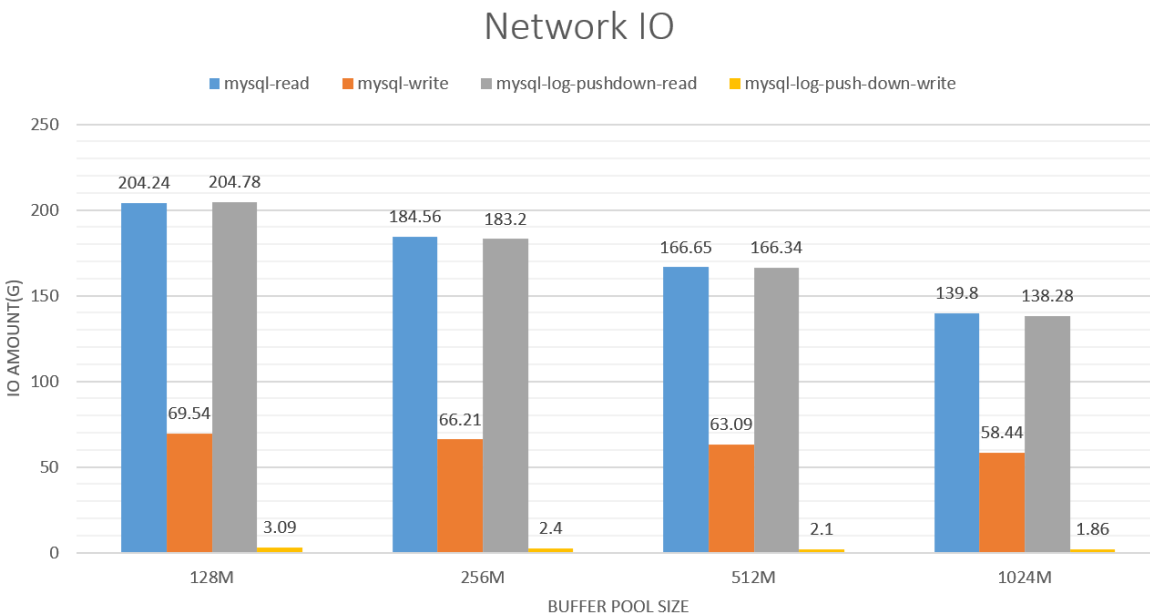
## 5.1Throughput test



picture8. TPC-cThroughput test

As shown above, in the fixed MySQL, the number of connections is 32. In the case of multiple threads, we will MySQL's Buffer Pool from 128M gradually increase to 1024M, the figure above shows TPC-C throughput comparison, green line represents the throughput in the original scenario, mysql-log-pushdown represents the throughput of the solution proposed in this paper. It can be seen that compared with the original scenario, the throughput is increased by nearly two times.

## 5.2network IO test



picture9.network IO test

As shown above, as shown, we modified TPC-C's source code in MySQL. The number of connections is 32. In the case of multiple threads, make each thread run 3000 transactions, after the test we measured the network between the computing nodes and the storage nodes. IOs shown in the figure above, compared with the original scene, the network IO basically remain unchanged, but because our solution basically filters out all data page write operations, only Redo Log, Undo Log and some other pages are written, so write IO is very small, compared to the original scene. Buffer Pool Size for 512M. In the case of IO, the most decreased 79.2 times. When running the same number of transactions, the larger the Buffer Pool's ability to cache more pages, triggering IO the number of times will also be less, so as Buffer Pool's size increases, the corresponding reading and writing IO the total amount shows a downward trend.