

# 分析mysql acid设计实现

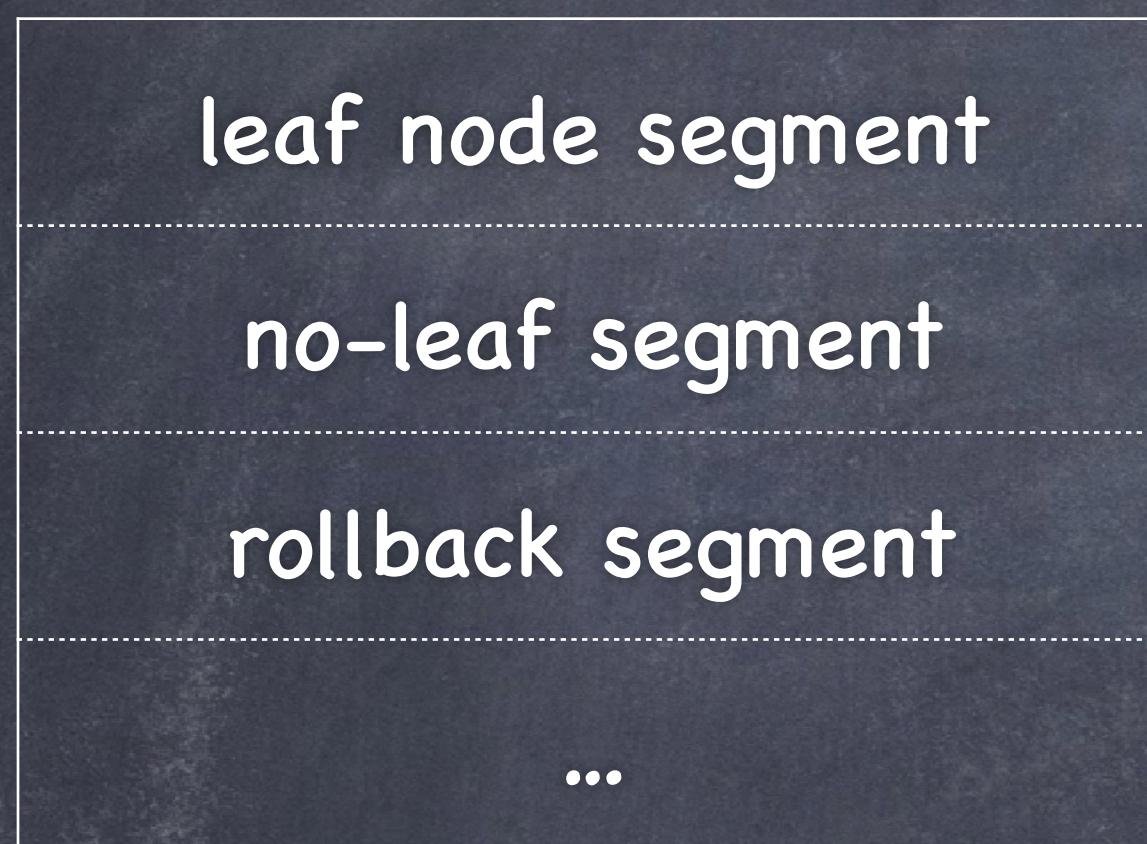
blog- [xiaorui.cc](http://xiaorui.cc)  
author- rfyiamcool

# List

- ② innodb 基本结构
- ② innodb acid 概念
- ② innodb buffer pool
- ② mvcc
- ② 各种延伸
- ② 锁, 各种锁

# innodb 的抽象结构

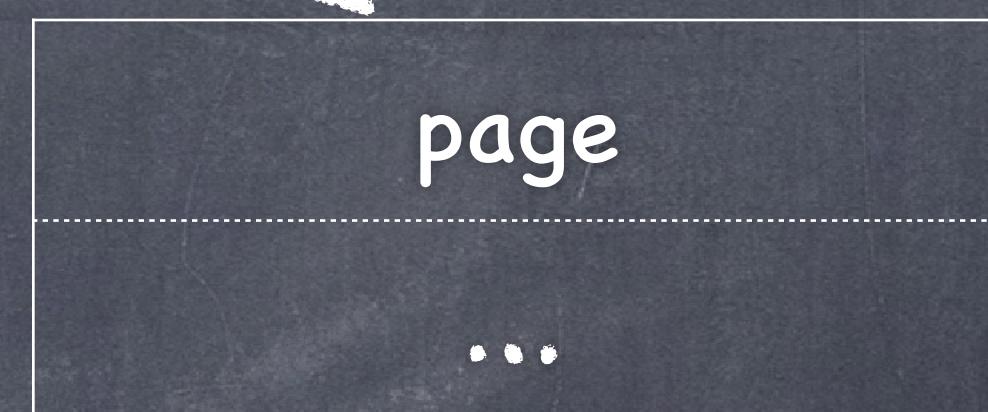
tablespace



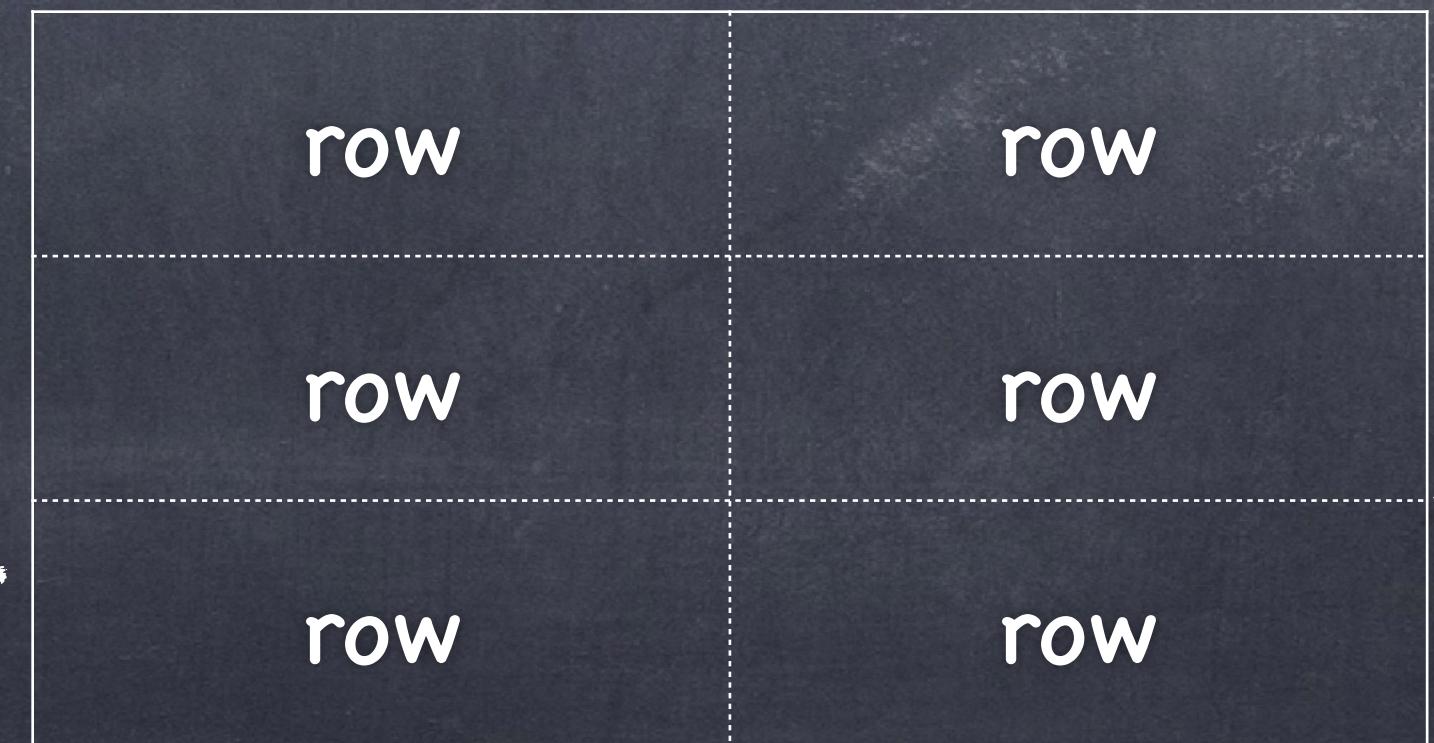
segment



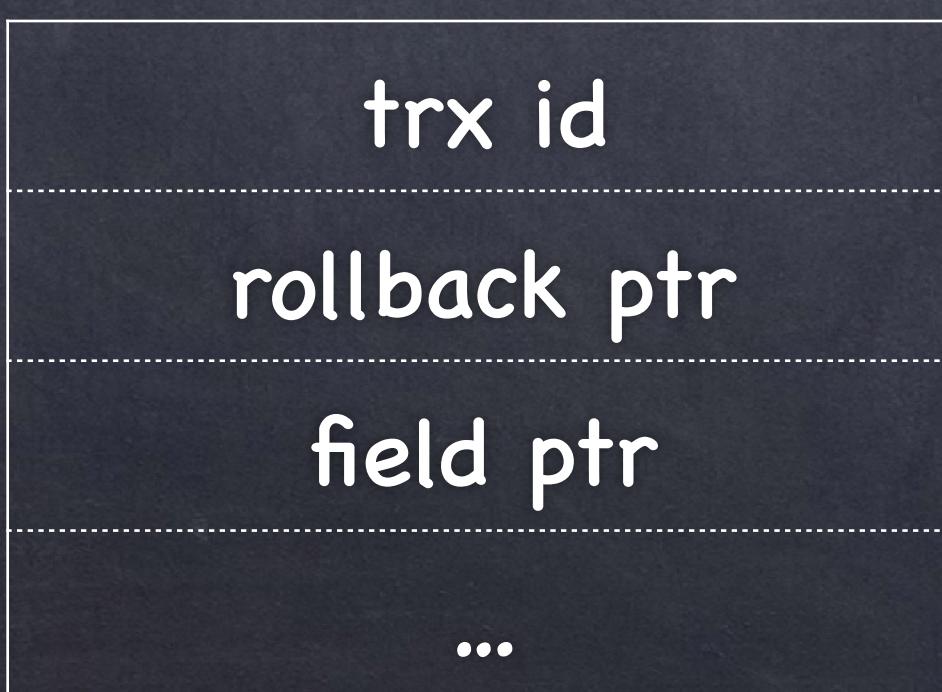
extent



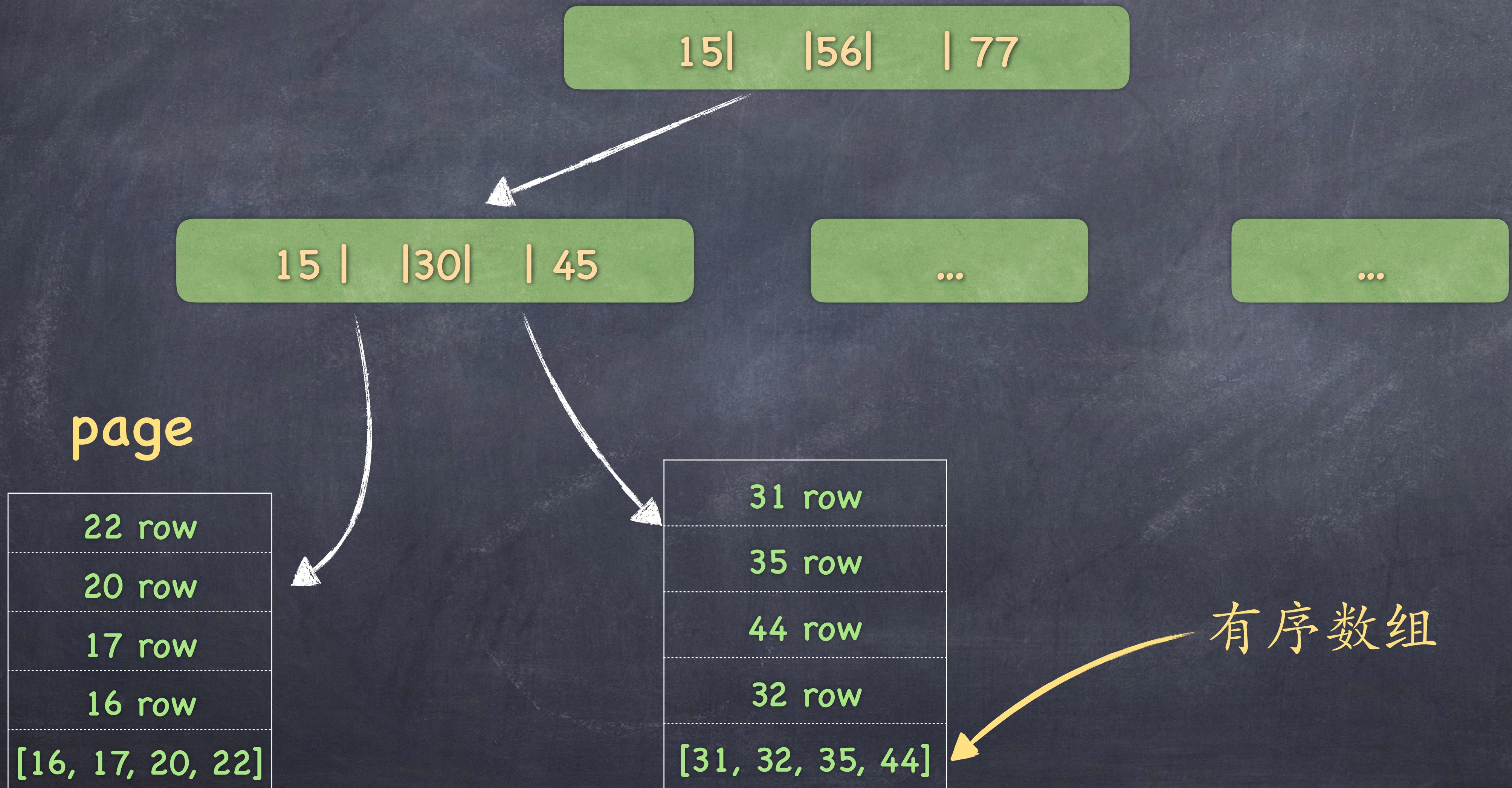
page



row



# b+tree 数据结构



# file system vs disk

innodb engine - - - > 16k

---

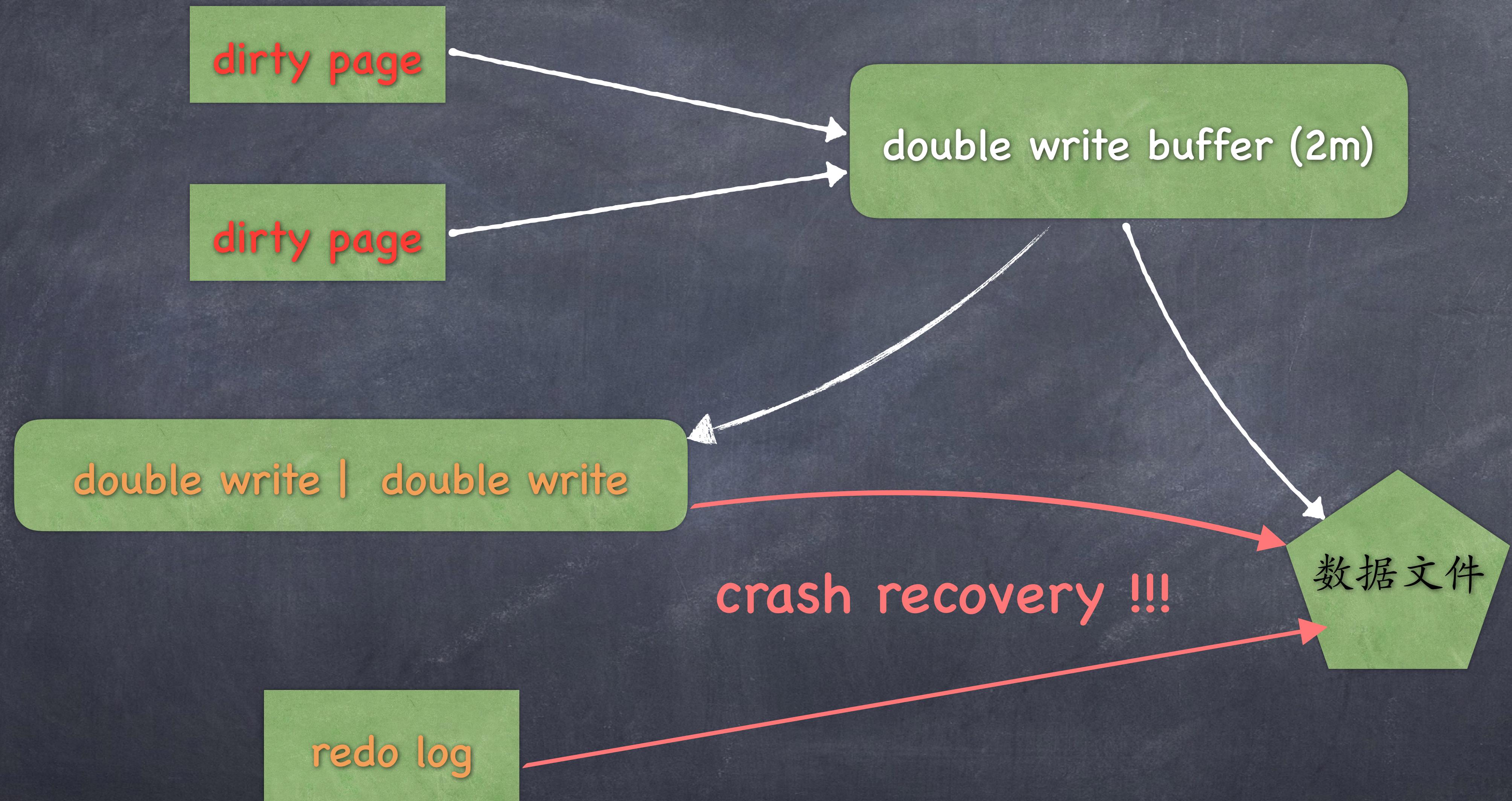
file system - - -> 4k

---

disk - - - > 512 bytes

# size 不同如何上下对齐？

- ④ partial page write ?
- ④ 解决办法是 doublewrite buffer + doublewrite
- ④ doublewrite 性能影响? 5% - 10%



# what is acid

- atomicity
- consistency
- isolation
- durability

# atomicity

- ◎ 一个事务为一个原子 不可分割
- ◎ 要么成功 要么失败

# consistency

- ② 经过一系列改动及异常崩溃，最后的结果是我要的.
- ② 最终一致性
- ② 强一致性

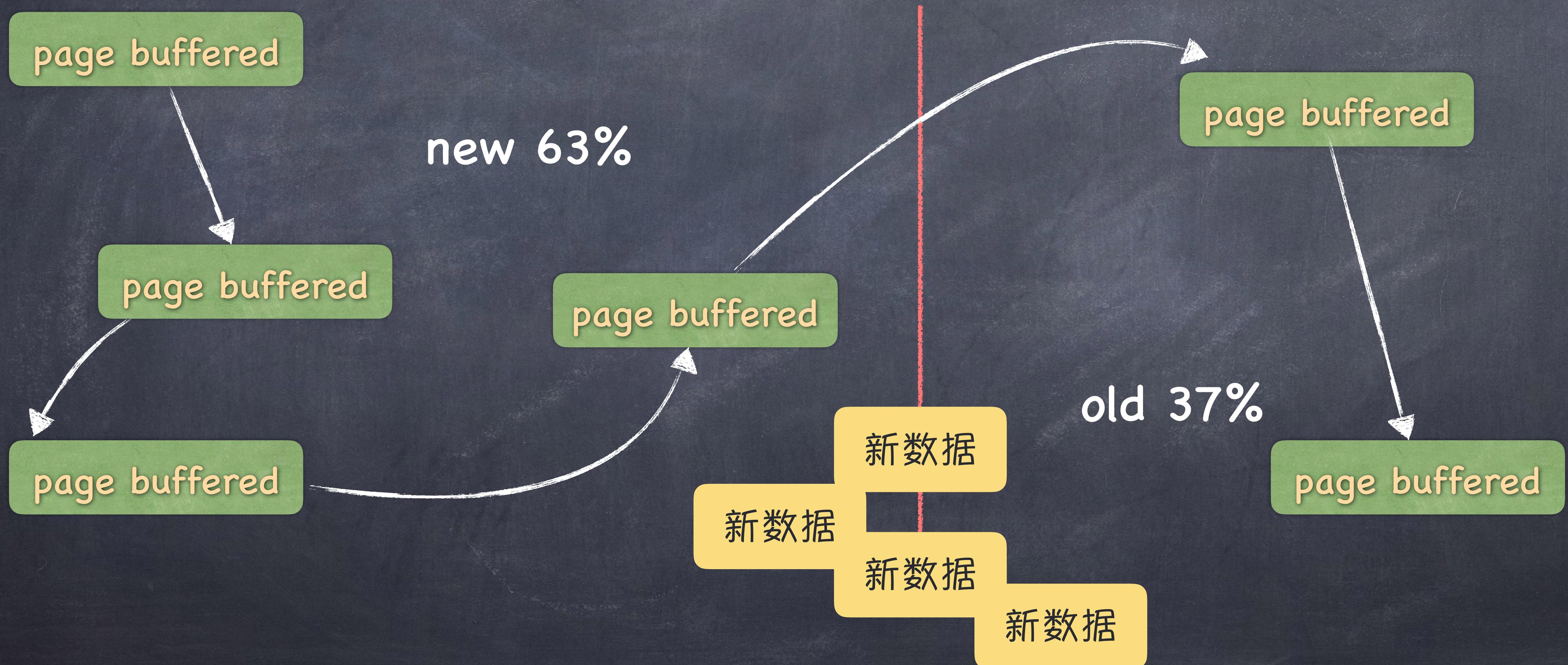
# isolation

- ⌚ mvcc
- ⌚ 隔离级别
  - ⌚ READ UNCOMMITTED
  - ⌚ READ COMMITTED
  - ⌚ REPEATABLE READ
  - ⌚ SERIALIZABLE

# durability

- 一旦提交成功，那么事务涉及的数据都会持久化
- 即使系统崩溃，修改的数据也不会丢失

# innodb buffer pool



update where blog = 'xiaorui.cc'

data in buffer pool

no !

Load page to buffer

yes !

update data and mark **dirty page !!!**

# innodb buffer pool

- ⦿ lru
- ⦿ dirty page
- ⦿ free list
- ⦿ flush list

# 一些概念总览

- ⌚ redo
- ⌚ undo
- ⌚ checkpoint
- ⌚ purge

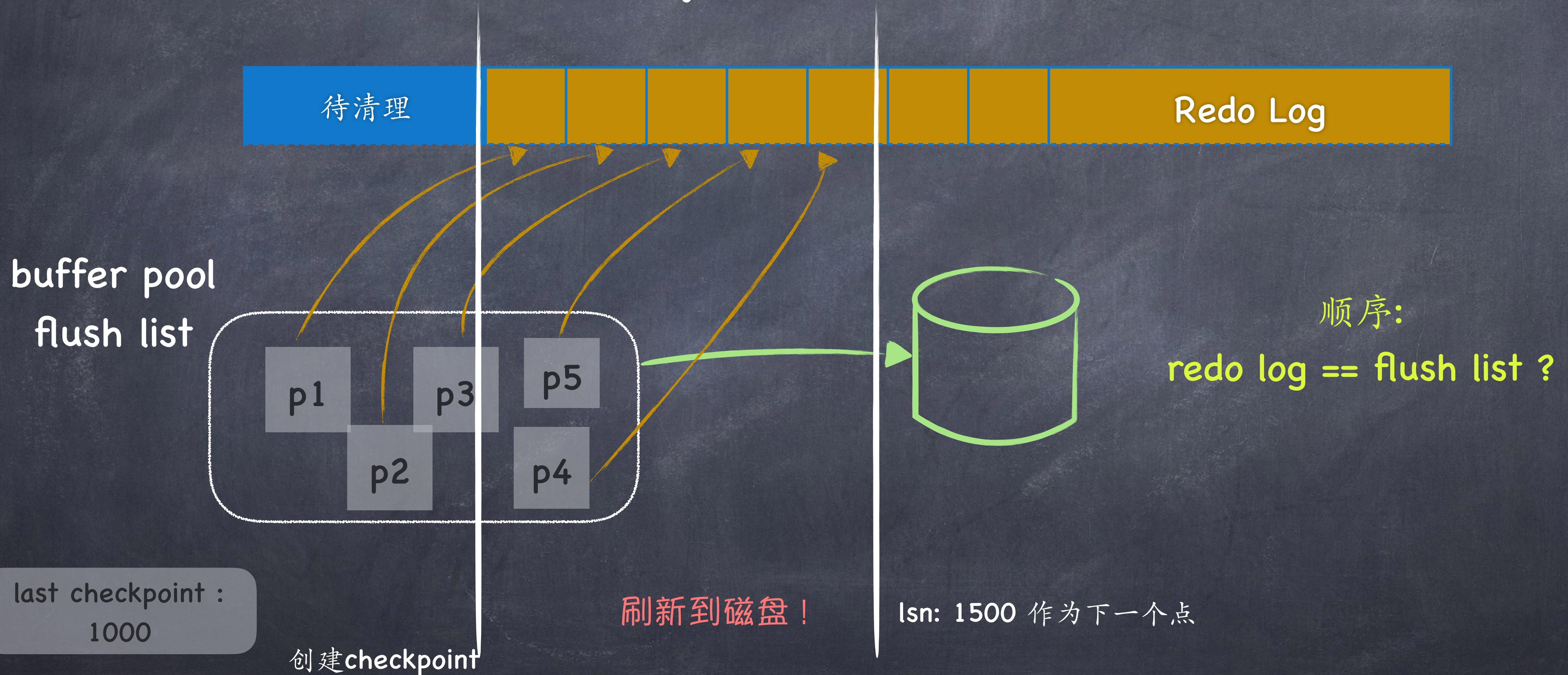
# if not checkpoint

- ② 写操作都是在buffer pool进行，每次提交都会顺序io增加redo日志.
- ③ buffer pool里面的dirty page越来越多 .
- ④ redo 越来越大，导致recovery时间过长 .

# if checkpoint

- ◎ 当 buffer pool 内存不够用时，将dirty page刷新到文件。
- ◎ 当redo文件过大又不可循环写入时，同上
- ◎ 缩短crash recovery时间

# checkpoint + lsn



# 一个事务过程

1. begin
2. undo log
3. update where id = xxx;
4. redo buffer
5. redo fsync; binlog; commit

if crash ?

if crash ?

if crash ?

one of three happen crash ?

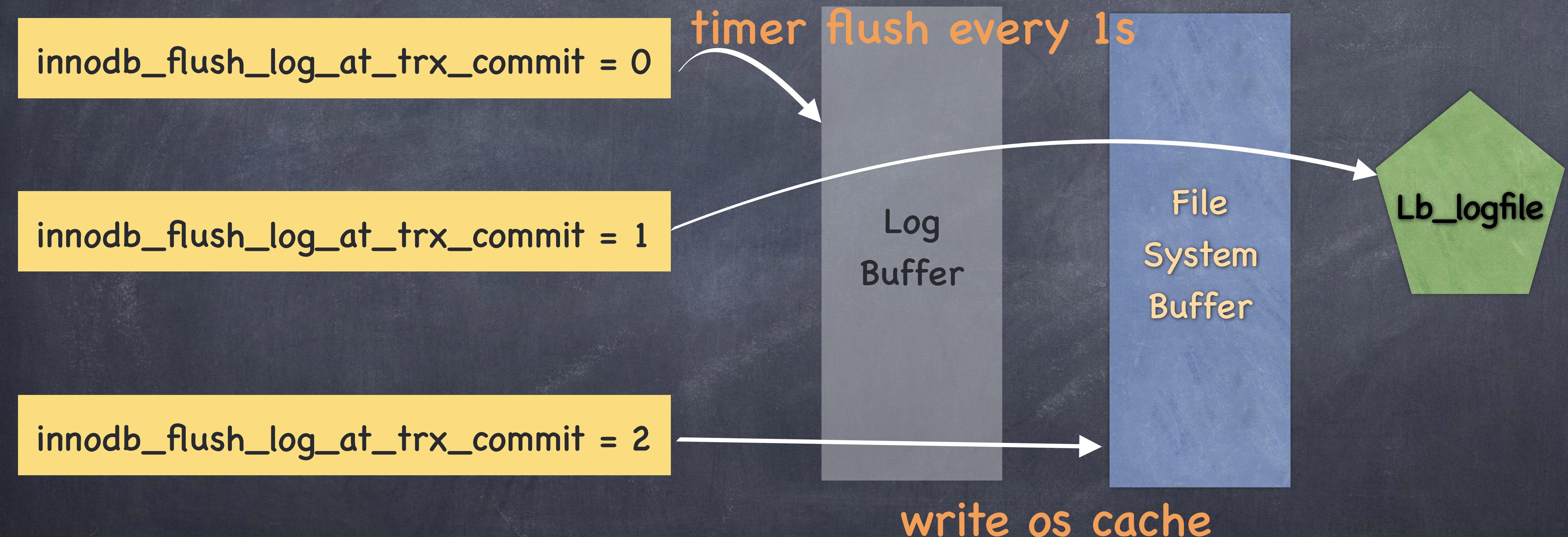
# redo

- ⌚ when redo save ?
- ⌚ when commit !
- ⌚ 顺序io
- ⌚ page number + block
- ⌚ when redo reduce ?
- ⌚ Ringbuffer full
- ⌚ flush timer run

# undo

- ④ 随机io
- ④ 是否存入文件, 取决于**buffer pool**

# flush log rule



# 数据库并发处理的演变

- ⌚ 基于表的读写锁
- ⌚ 基于索引的读写锁（分离锁）
- ⌚ MVCC

mvcc

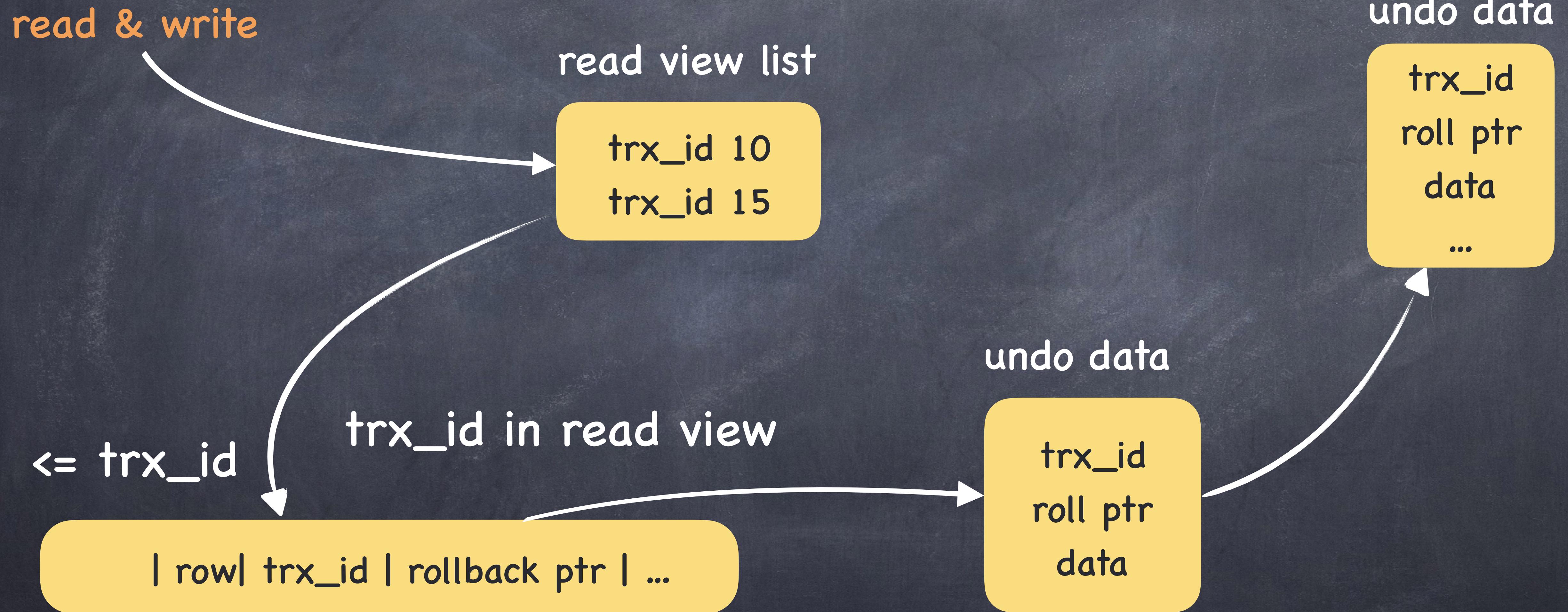
- ② 读读非阻塞
- ② 读写非阻塞
- ② 写写阻塞

一句话, 并发的效果, 串行的结果 !!!

# 又是概念

- ⑥ read view
- ⑥ 隐藏字段
- ⑥ 是否可见的条件
- ⑥ 当前读
- ⑥ 快照读

# mvcc 流程图



# 幻读

- ⌚ record lock
- ⌚ gap lock
- ⌚ next gap lock

# 业务数据不一致

事务 A

```
ts begin  
ts select stock where ticket=10 stock=1  
ts update stock=stock-1 stock=0  
ts commit  
ts  
ts  
ts
```

事务 B

```
begin  
select stock where ticket = 10  
update stock=stock-1 stock=-1  
commit
```

最后 -1

# 互斥锁 vs 乐观锁

- ⌚ 互斥锁
  - ⌚ `select * from train where ticket = 10 for update;`
- ⌚ 乐观锁
  - ⌚ `while 1:`
  - ⌚ `update train set stock=stock-1 where ticket = 10 and stock = x;`

# crash recovery

- ④ 从**doublewrite** 修正缺斤少两的**page**
- ④ 重做提交的**redo** 日志
- ④ 回滚执行一半的 未提交的 事务

“Q & A”

-fengyun rui