

# 大数处理

## 常见海量处理题目解题关键

- 1、分而治之。通过哈希函数将大任务分流到机器，或分流成小文件。
- 2、常用的hashMap或bitmap。

难点：通讯、时间和空间的估算。

请对10亿个IPV4的ip地址进行排序，每个ip只会出现一次。

IPV4的ip数量 $\approx 42$ 亿

申请长度为 $2^{32}$ 的bit类型的数组。



(每个位置上是一个bit, 只可表示0或1两种状态)  
长度为 $2^{32}$ 的bit数组, 空间约为128m。

请对10亿人的年龄进行排序。



有一个包含20亿个全是32位整数的大文件，在其中找到出现次数最多的数。但是内存限制只有2G。

## 20亿个32位整数的大文件

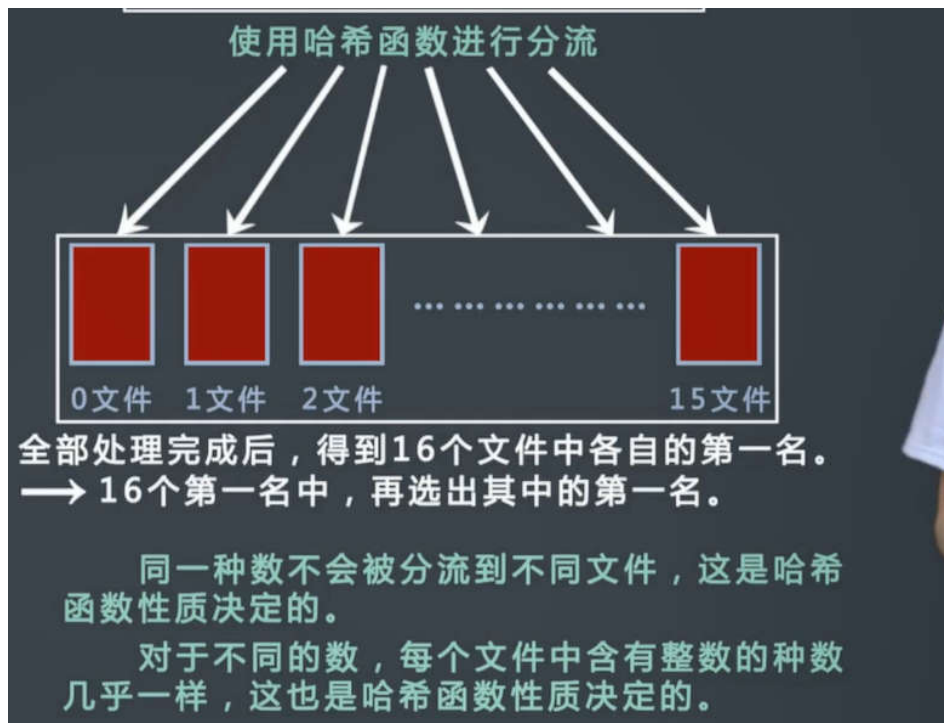
hashmap记录所有数出现的次数。

key → 具体某一种数  
value → 这种数出现的次数

hashmap记录所有数出现的次数。

4字节 ← 整型 ← key → 具体某一种数  
4字节 ← 整型 ← value → 这种数出现的次数  
一条记录 ( key , value ) 占有8字节  
记录条数为2亿时，大约1.6G内存。

所以用哈希表来直接统计20亿个整数的方案，  
会导致内存不足。

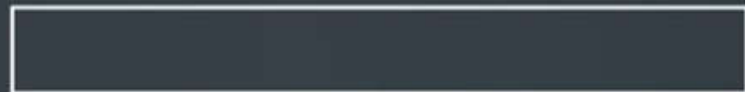


32位无符号整数的范围是0~4294967295。现在有一个正好包含40亿个无符号整数的文件，所以在整个范围中必然有没出现过的数。可以使用最多10M的内存，只用找到一个没出现过的数即可，该如何找？

如果用哈希表来记录所有的数，最差情况下，  
将出现40亿个不同的数。

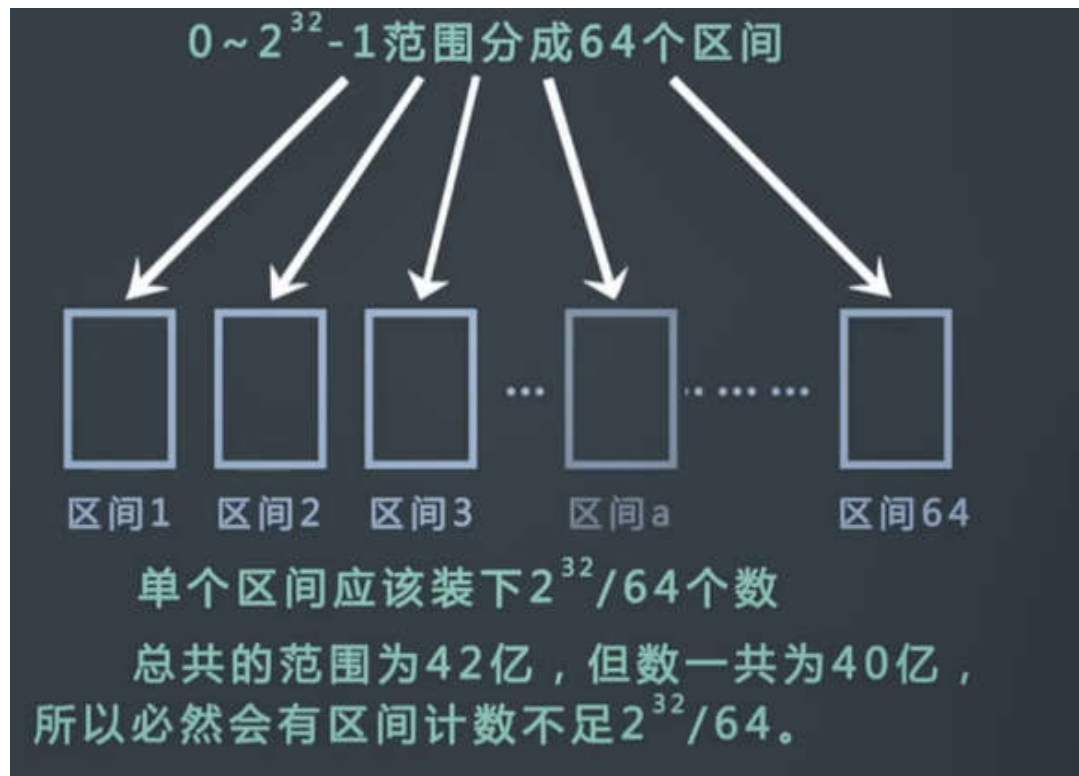
每一条记录占有4字节，大约需要16G内存。

bitmap:



0 1 2 ... ..  $2^{32}-1$

每个位置为1个bit，只能表示0或1两种状态。  
大约占用500m空间。





区间a

再遍历一次40个数，此时只关注区间a上的数，并用bitmap统计区间a上的数的出现情况。

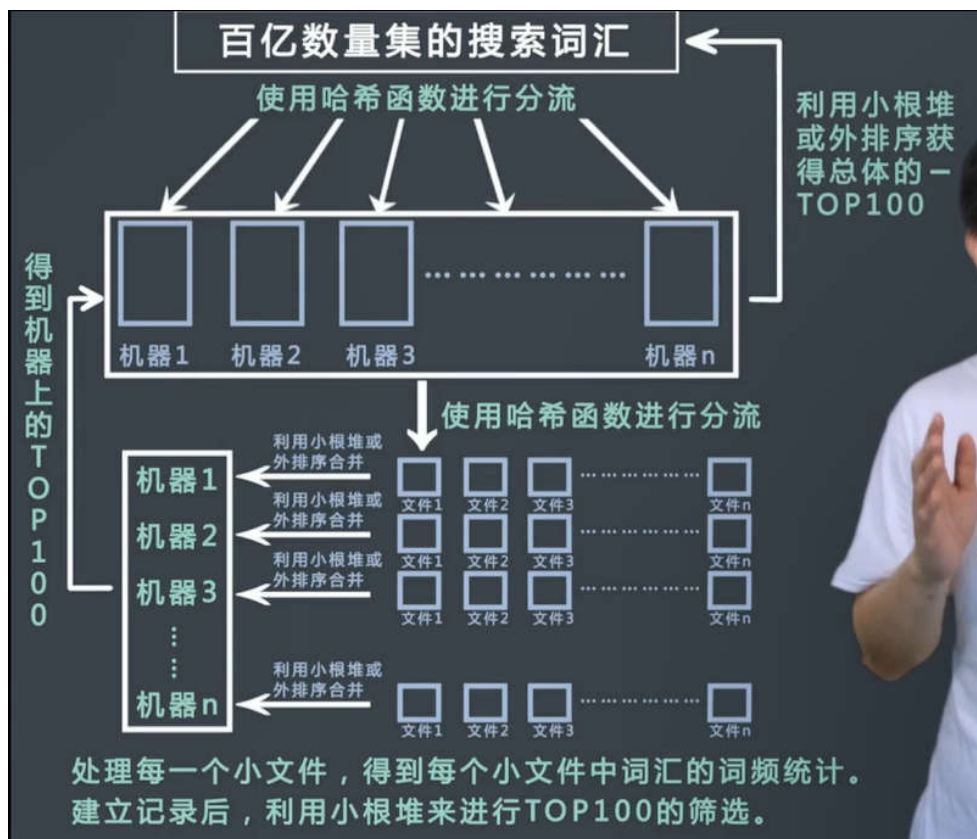
占用差不多8m空间。

## 总结：

- 1、根据内存限制决定区间大小，根据区间大小，得到有多少个变量，来记录每个区间的数出现的次数。
- 2、统计区间上的数的出现次数，找到不足的区间。
- 3、利用bitmap对不满的区间，进行这个区间上的数的词频统计。



某搜索公司一天的用户搜索词汇是海量的，假设有百亿的数据量，请设计一种求出每天最热100词的可行办法。



工程师常使用服务器集群来设计和实现数据缓存，以下是常见的策略。1，无论是添加、查询还是删除数据，都先将数据的id通过哈希函数转换成一个哈希值，记为key。2，如果目前机器有N台，则计算 $\text{key} \% N$ 的值，这个值就是该数据所属的机器编号，无论是添加、删除还是查询操作，都只在这台机器上进行。请分析这种缓存策略可能带来的问题，并提出改进的方案。

## 一致性哈希算法

