

本节内容

# 外部排序

王道考研/CSKAOYAN.COM

## 知识总览

外部排序

外存与内存之间的数据交换

外部排序的原理

影响外部排序效率的因素

优化思路

王道考研/CSKAOYAN.COM

## 外存、内存之间的数据交换

操作系统以“块”为单位对磁盘存储空间进行管理，如：每块大小 1KB  
各个磁盘块内存放着各种各样的数据



1KB缓冲区

内存



王道考研/CSKAOYAN.COM

## 外存、内存之间的数据交换

操作系统以“块”为单位对磁盘存储空间进行管理，如：每块大小 1KB  
各个磁盘块内存放着各种各样的数据



1 2 3

内存



磁盘的读/写以“块”为单位  
数据读入内存后才能被修改  
修改完了还要写回磁盘

王道考研/CSKAOYAN.COM

## 外存、内存之间的数据交换

操作系统以“块”为单位对磁盘存储空间进行管理，如：每块大小 1KB  
各个磁盘块内存放着各种各样的数据



读磁盘

45 | 23 | 11

内存

写磁盘

磁盘的读/写以“块”为单位  
数据读入内存后才能被修改  
修改完了还要写回磁盘

王道考研/CSKAOYAN.COM

## 外存、内存之间的数据交换

操作系统以“块”为单位对磁盘存储空间进行管理，如：每块大小 1KB  
各个磁盘块内存放着各种各样的数据



读磁盘

45 | 23 | 11

内存

写磁盘

磁盘的读/写以“块”为单位  
数据读入内存后才能被修改  
修改完了还要写回磁盘

王道考研/CSKAOYAN.COM

## 外存、内存之间的数据交换

操作系统以“块”为单位对磁盘存储空间进行管理，如：每块大小 1KB  
各个磁盘块内存放着各种各样的数据



45 23 11

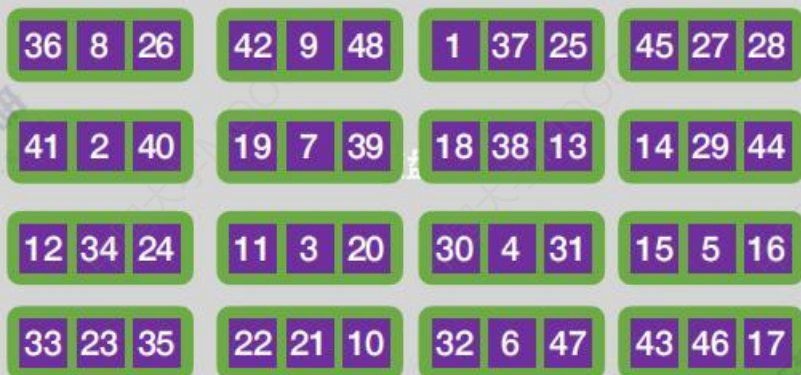
内存

磁盘的读/写以“块”为单位  
数据读入内存后才能被修改  
修改完了还要写回磁盘

王道考研/CSKAOYAN.COM

## 外部排序原理

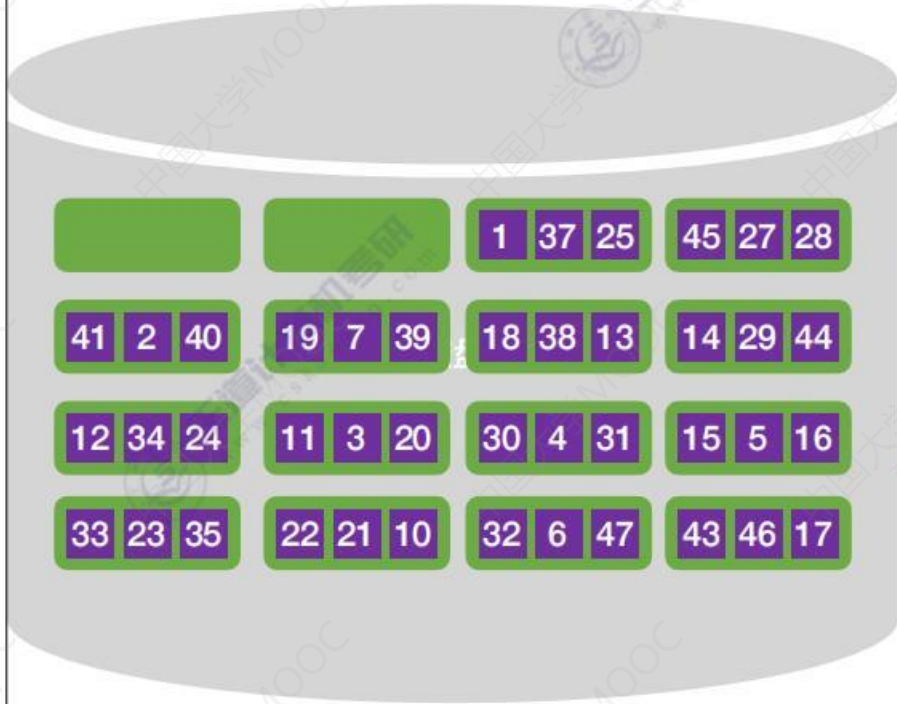
外部排序：数据元素太多，无法一次全部读入内存进行排序



使用“归并排序”的方法，最少只需在内存中分配3块大小的缓冲区即可对任意一个大文件进行排序

王道考研/CSKAOYAN.COM

### 构造初始“归并段”

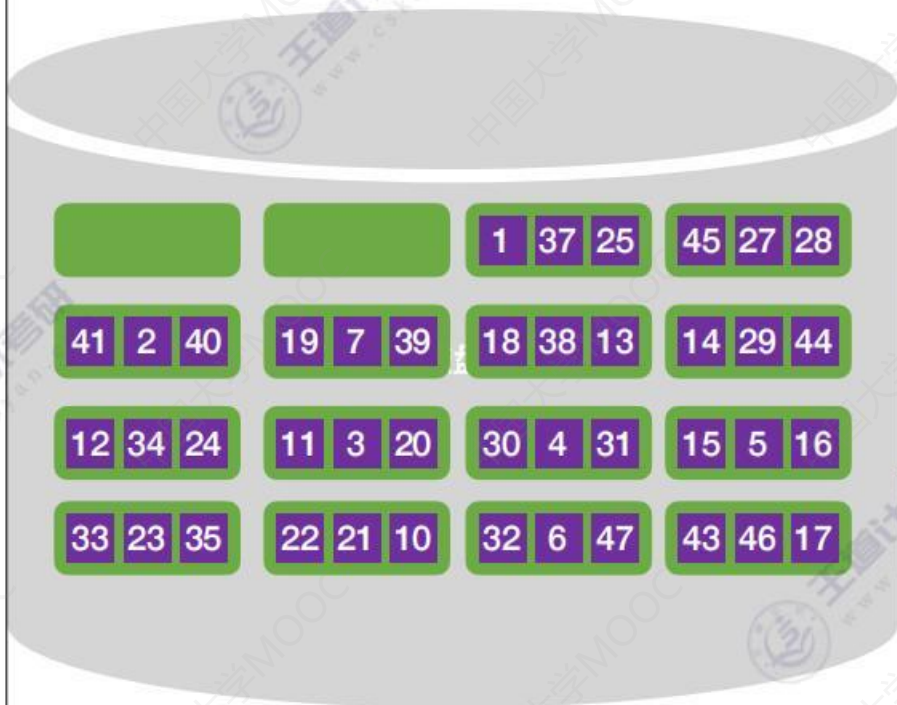


“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



王道考研/CSKAOYAN.COM

### 构造初始“归并段”

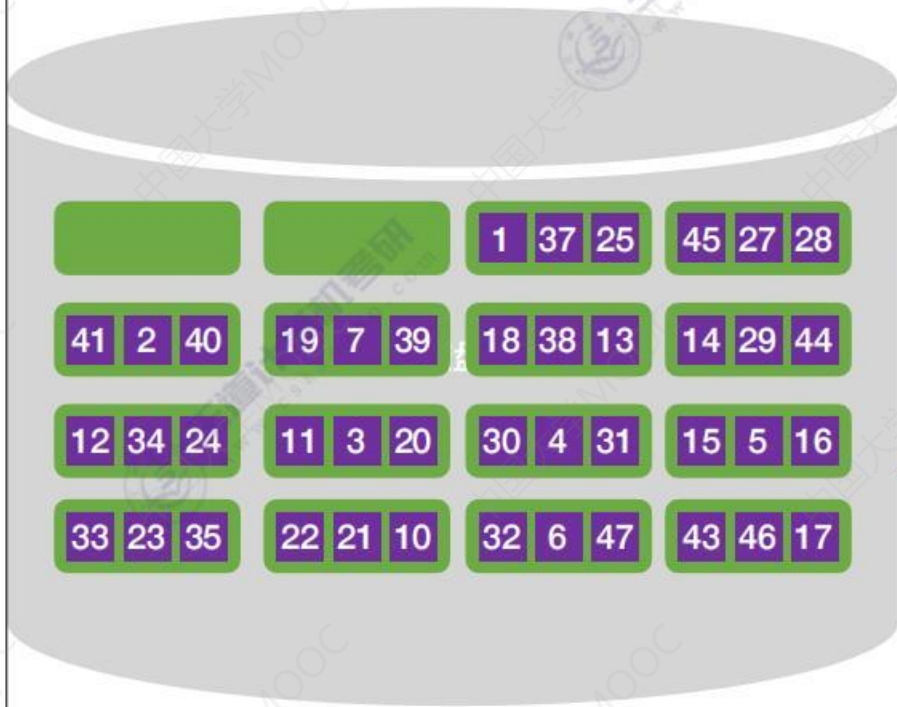


“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



王道考研/CSKAOYAN.COM

### 构造初始“归并段”

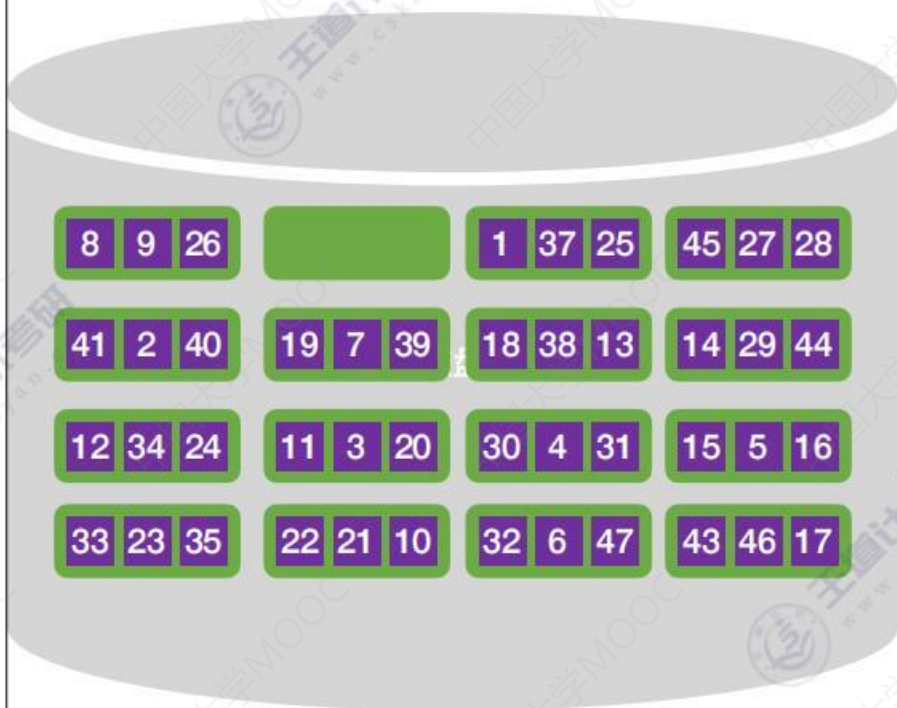


“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



王道考研/CSKAOYAN.COM

### 构造初始“归并段”

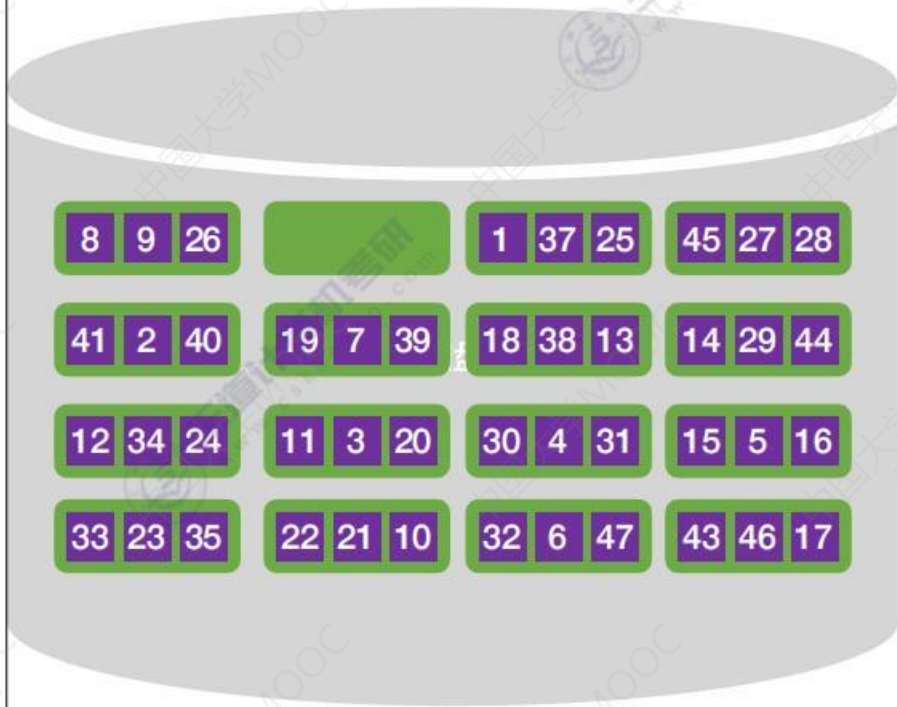


“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



王道考研/CSKAOYAN.COM

### 构造初始“归并段”

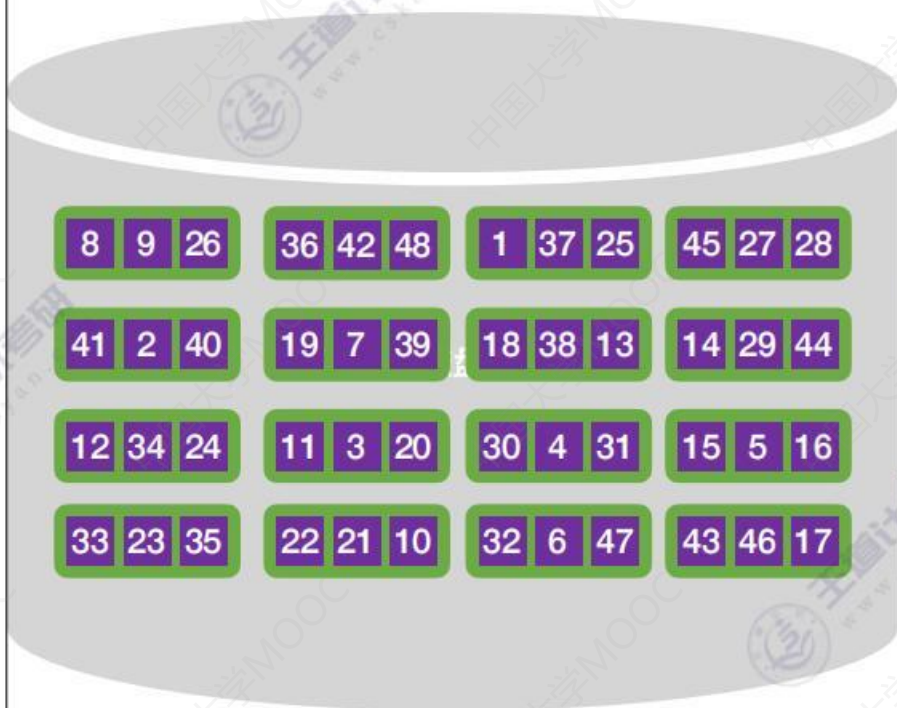


“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



王道考研/CSKAOYAN.COM

### 构造初始“归并段”



“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



王道考研/CSKAOYAN.COM

### 构造初始“归并段”

一个有序的“归并段”

8	9	26	36	42	48	1	37	25	45	27	28
41	2	40	19	7	39	18	38	13	14	29	44
12	34	24	11	3	20	30	4	31	15	5	16
33	23	35	22	21	10	32	6	47	43	46	17

读磁盘

写磁盘

“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘

输出缓冲区

输入缓冲区1

内存

输入缓冲区2

王道考研/CSKAOYAN.COM

### 构造初始“归并段”

8	9	26	36	42	48						
41	2	40	19	7	39	18	38	13	14	29	44
12	34	24	11	3	20	30	4	31	15	5	16
33	23	35	22	21	10	32	6	47	43	46	17

读磁盘

写磁盘

“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘

输出缓冲区

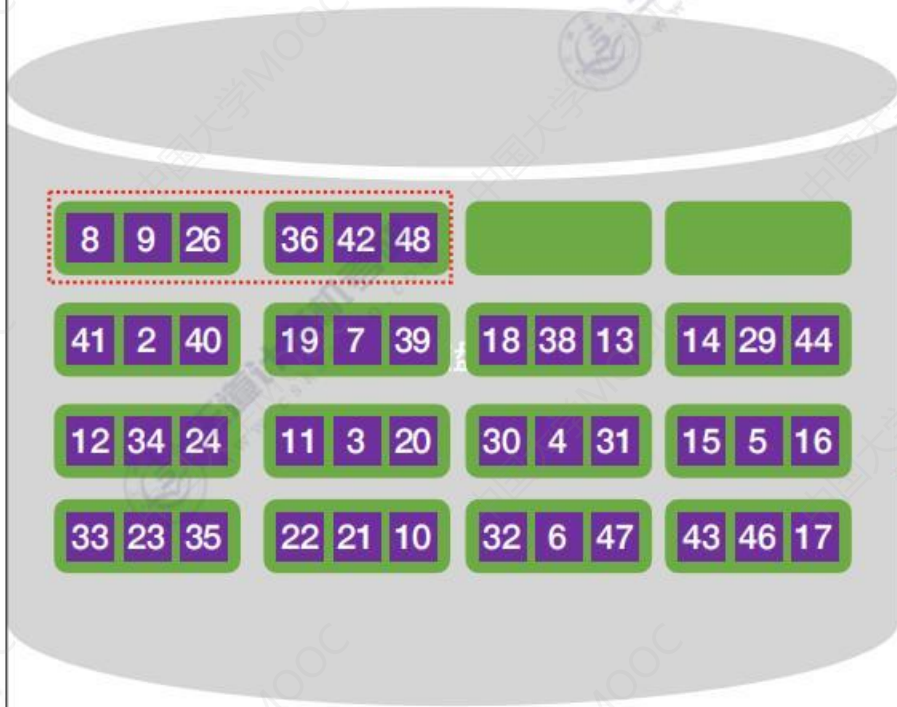
输入缓冲区1

内存

输入缓冲区2

王道考研/CSKAOYAN.COM

### 构造初始“归并段”

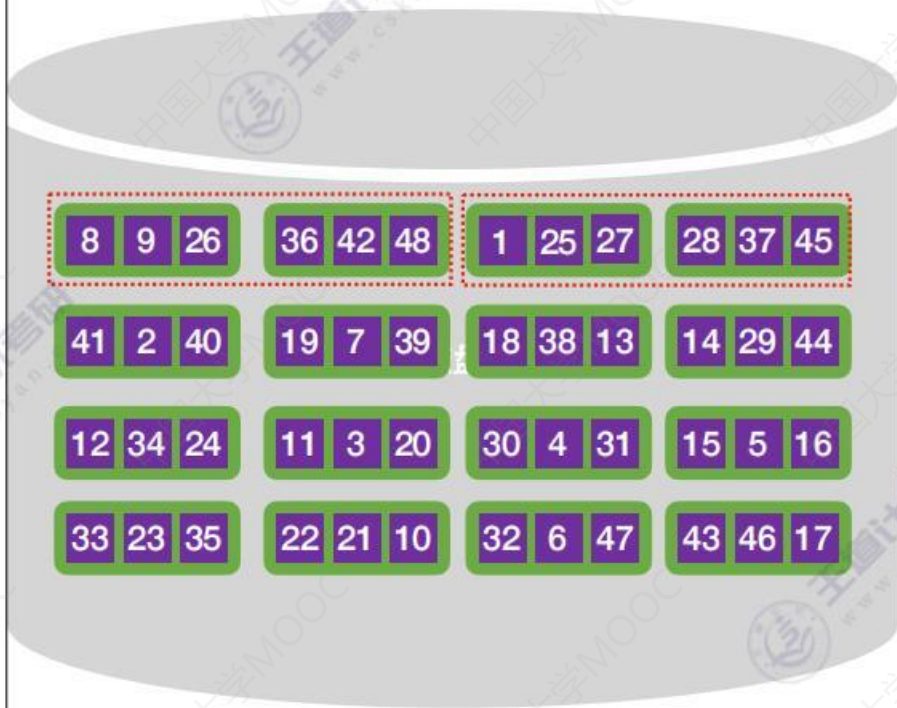


“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



王道考研/CSKAOYAN.COM

### 构造初始“归并段”



“归并排序”要求各个子序列有序，每次读入两个块的内容，进行内部排序后写回磁盘



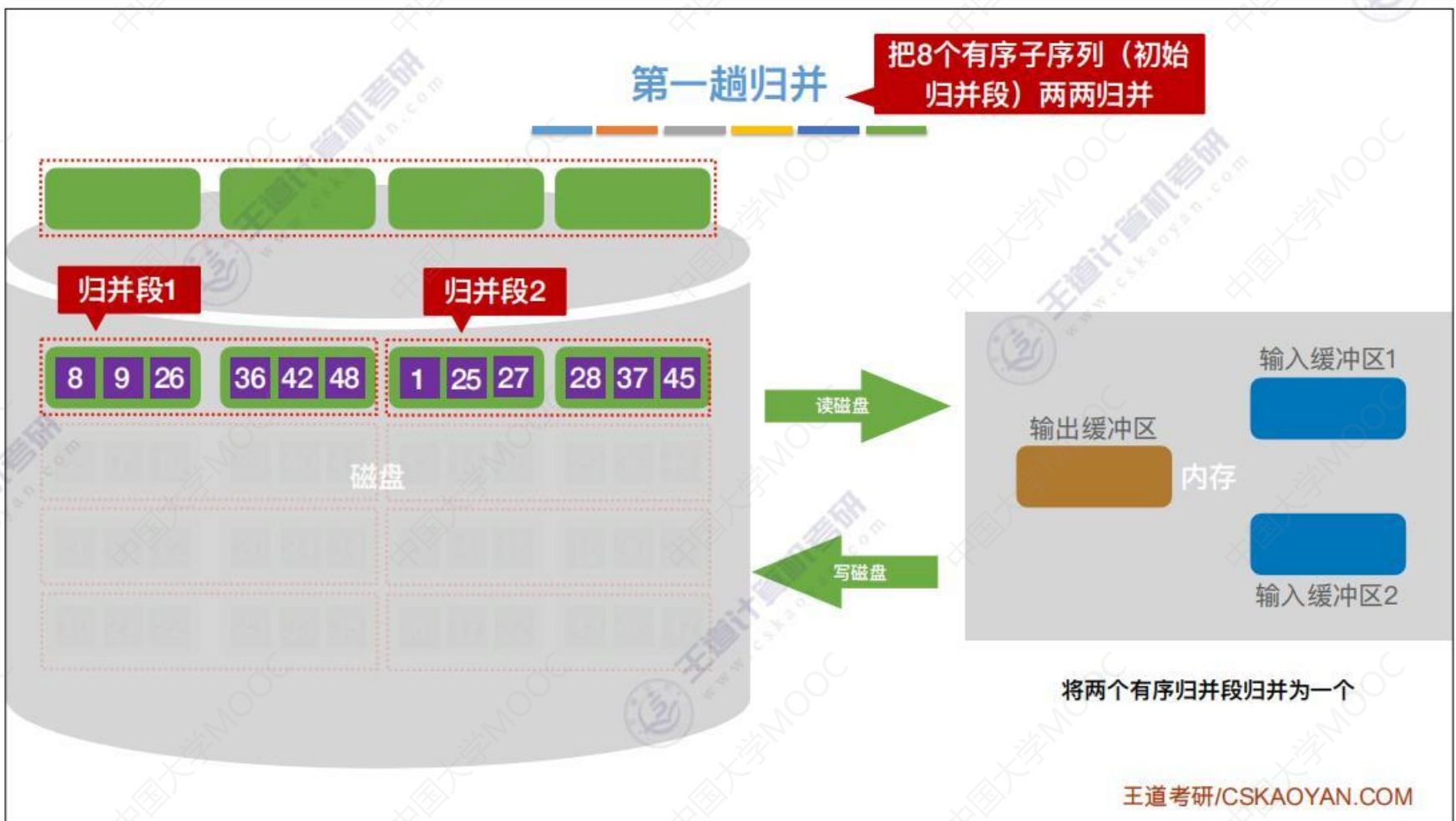
王道考研/CSKAOYAN.COM

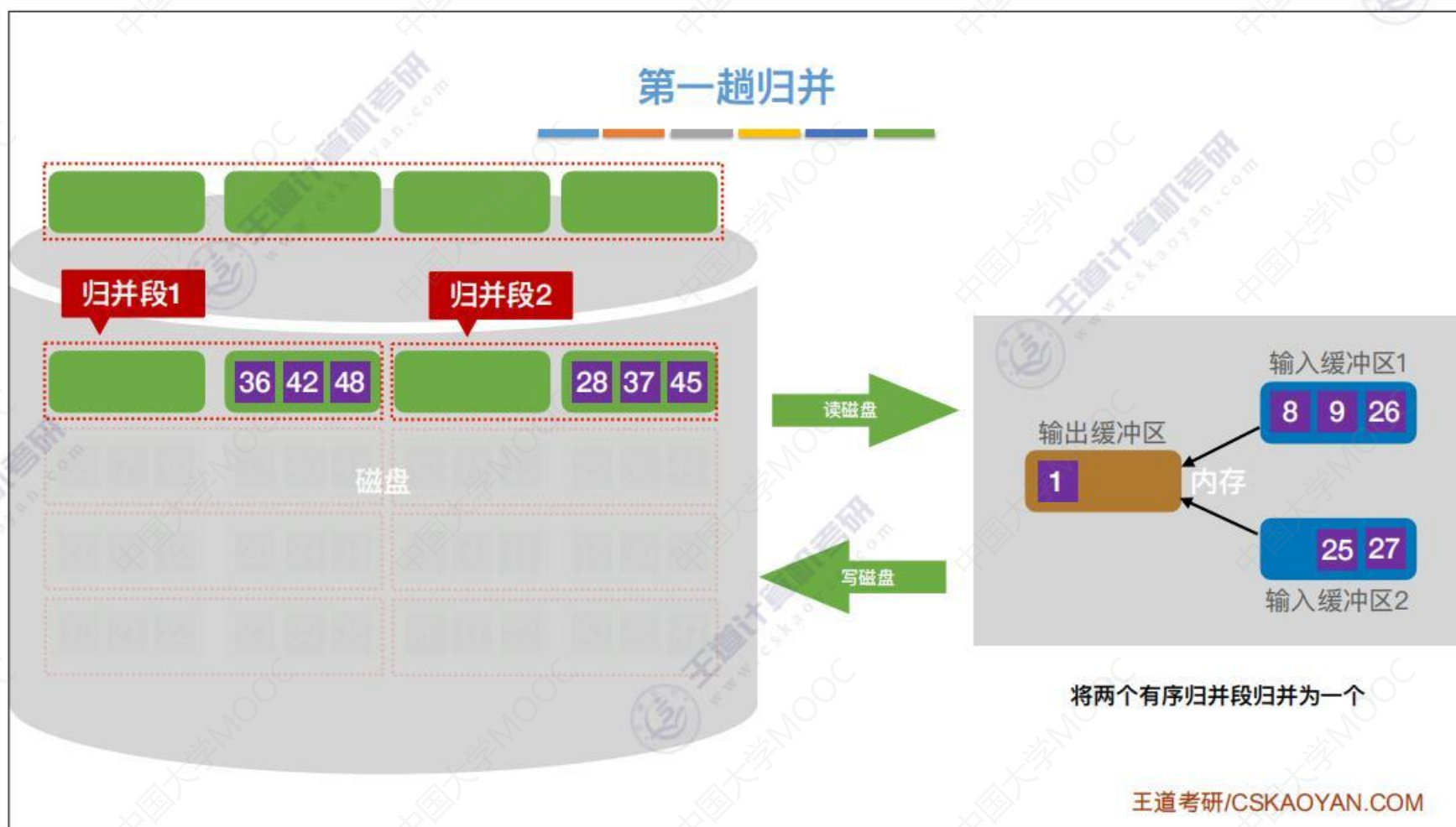
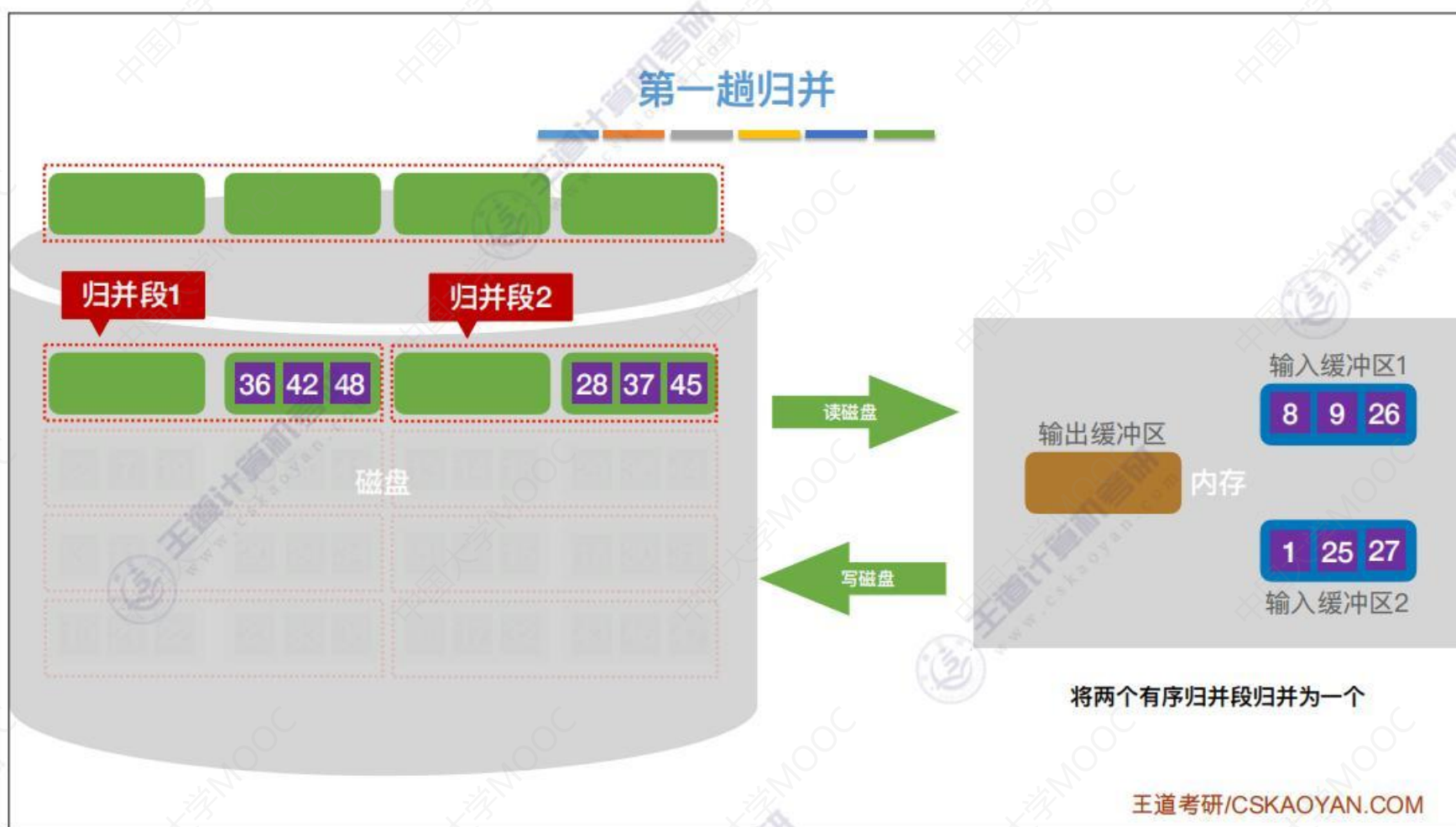
### 构造初始“归并段”

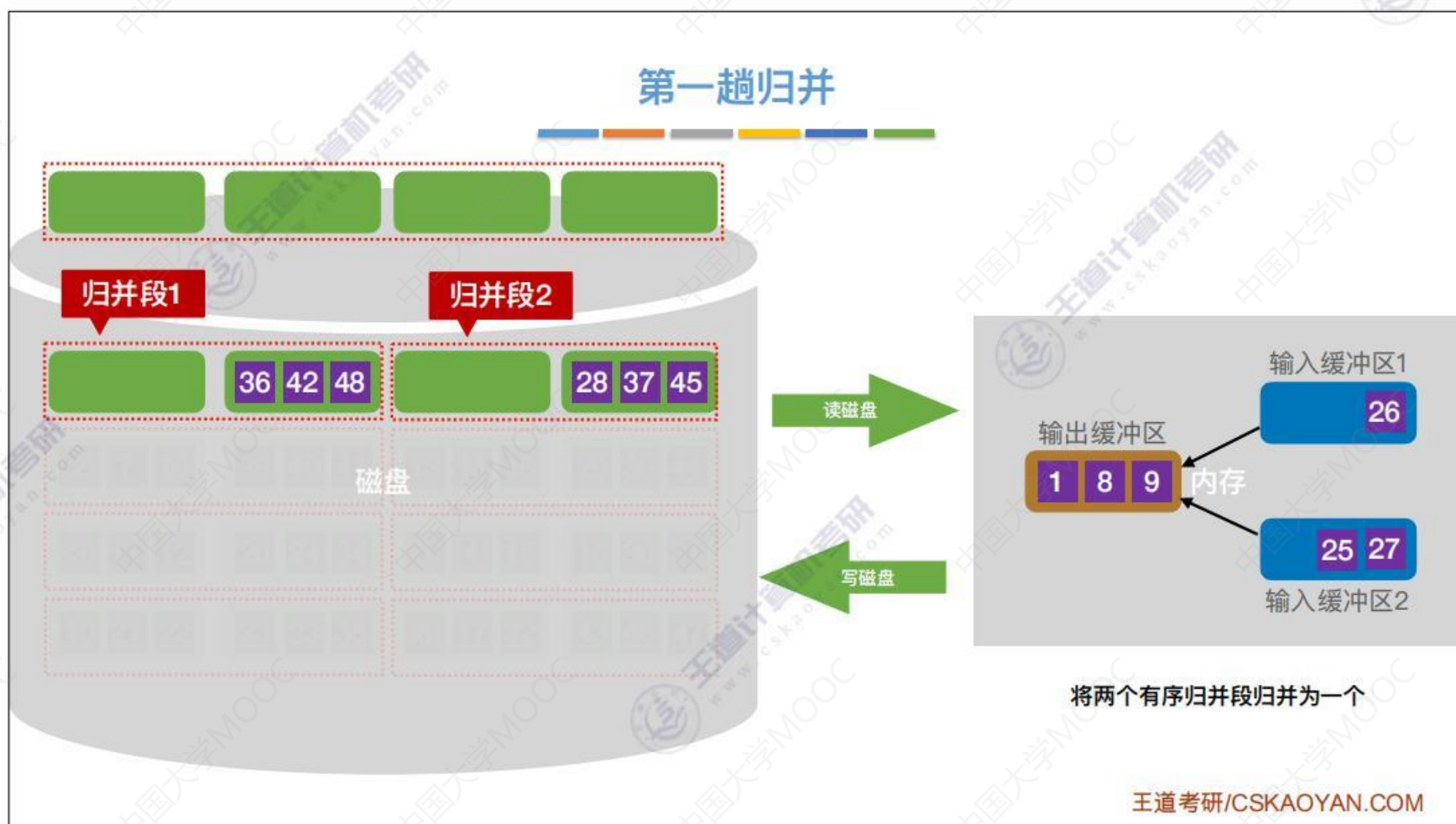
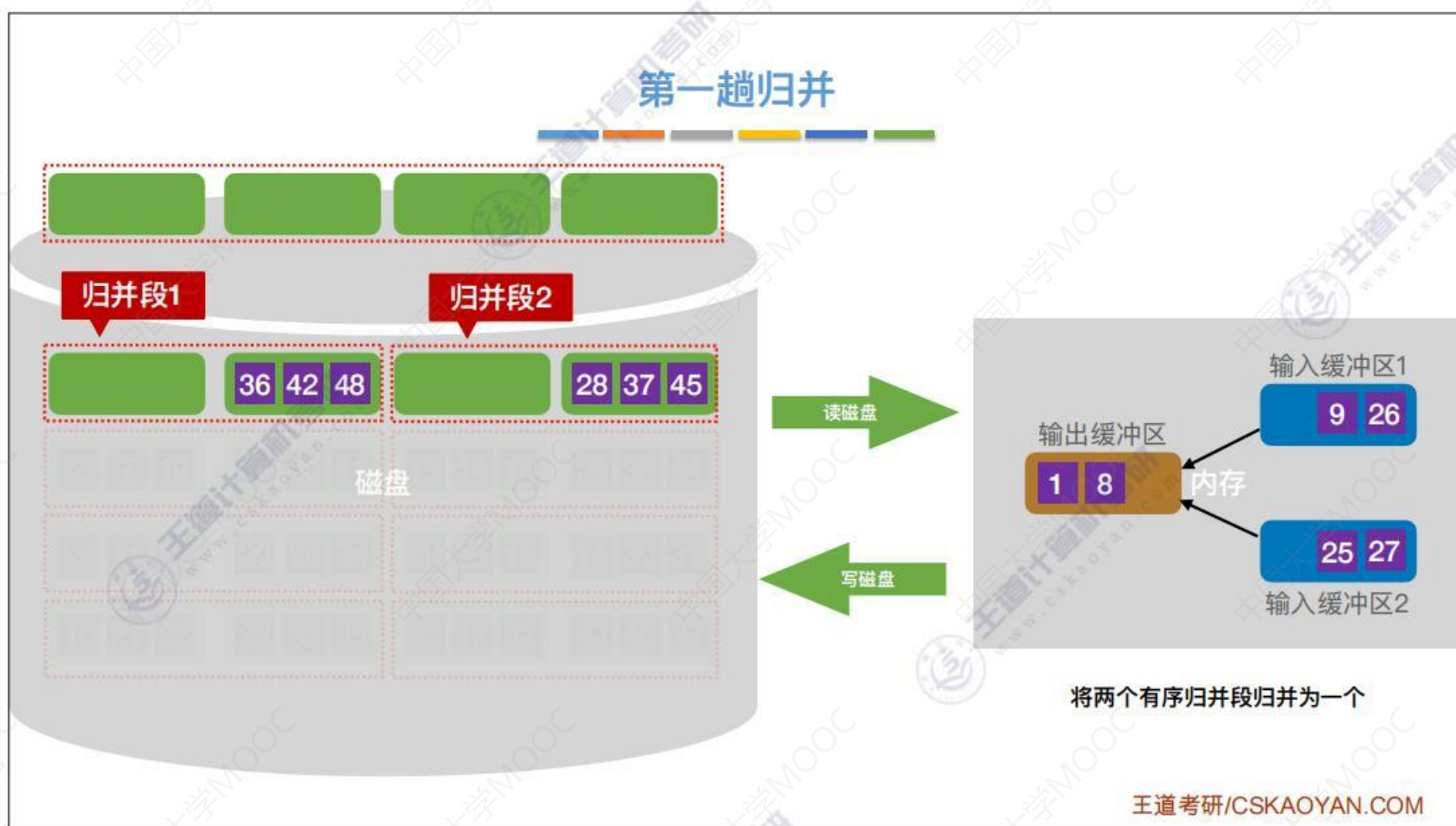


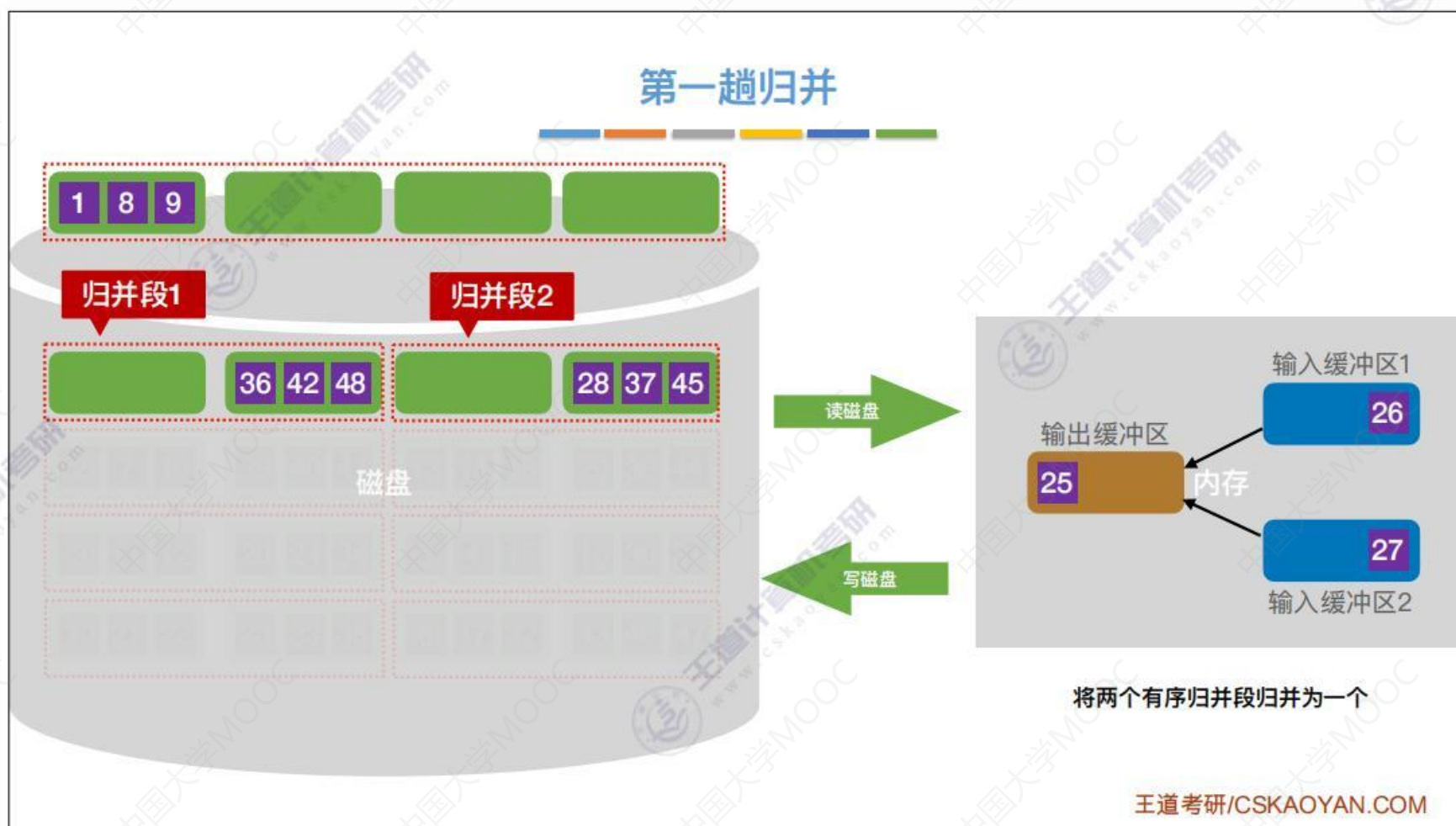
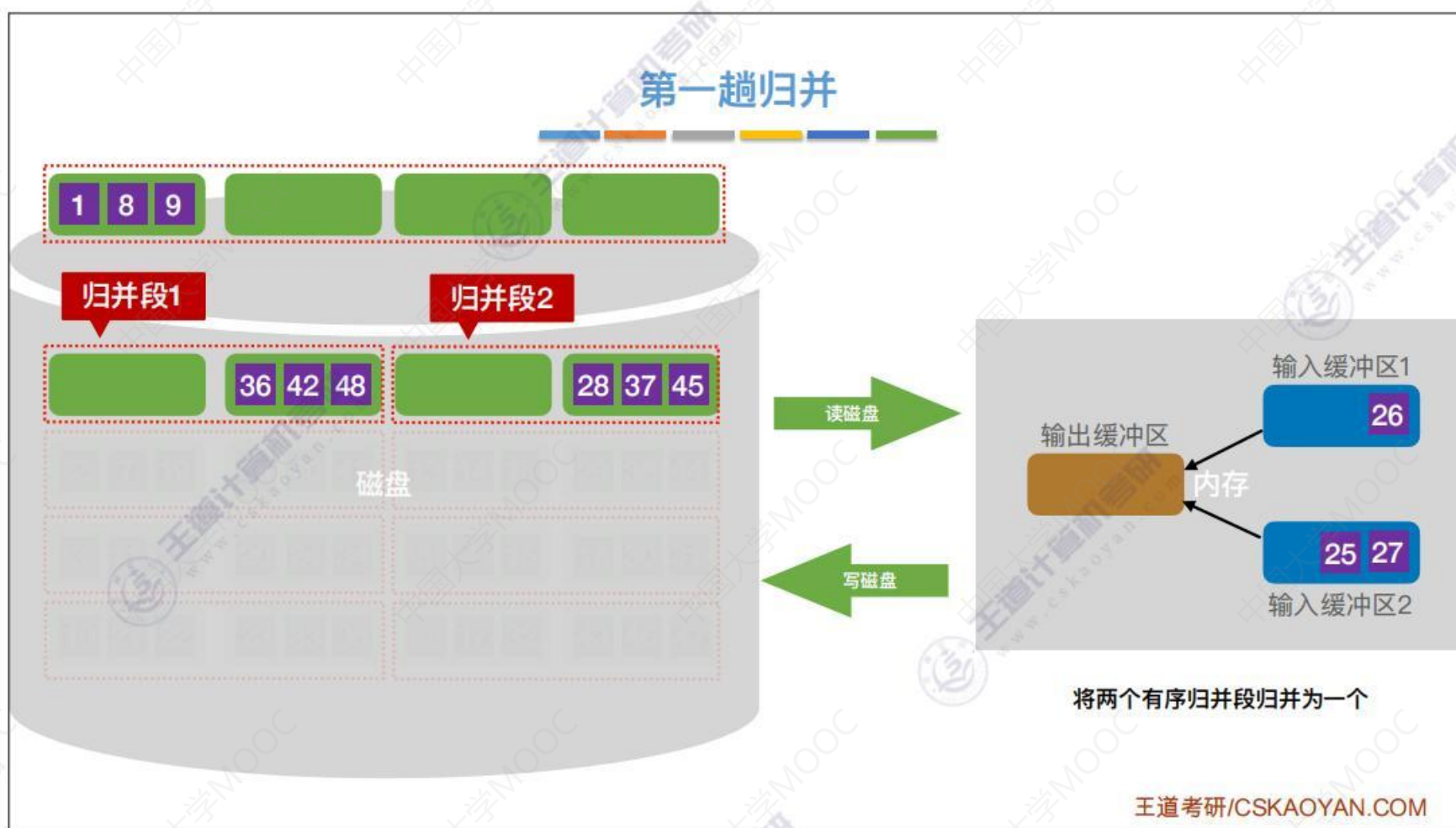
### 第一趟归并

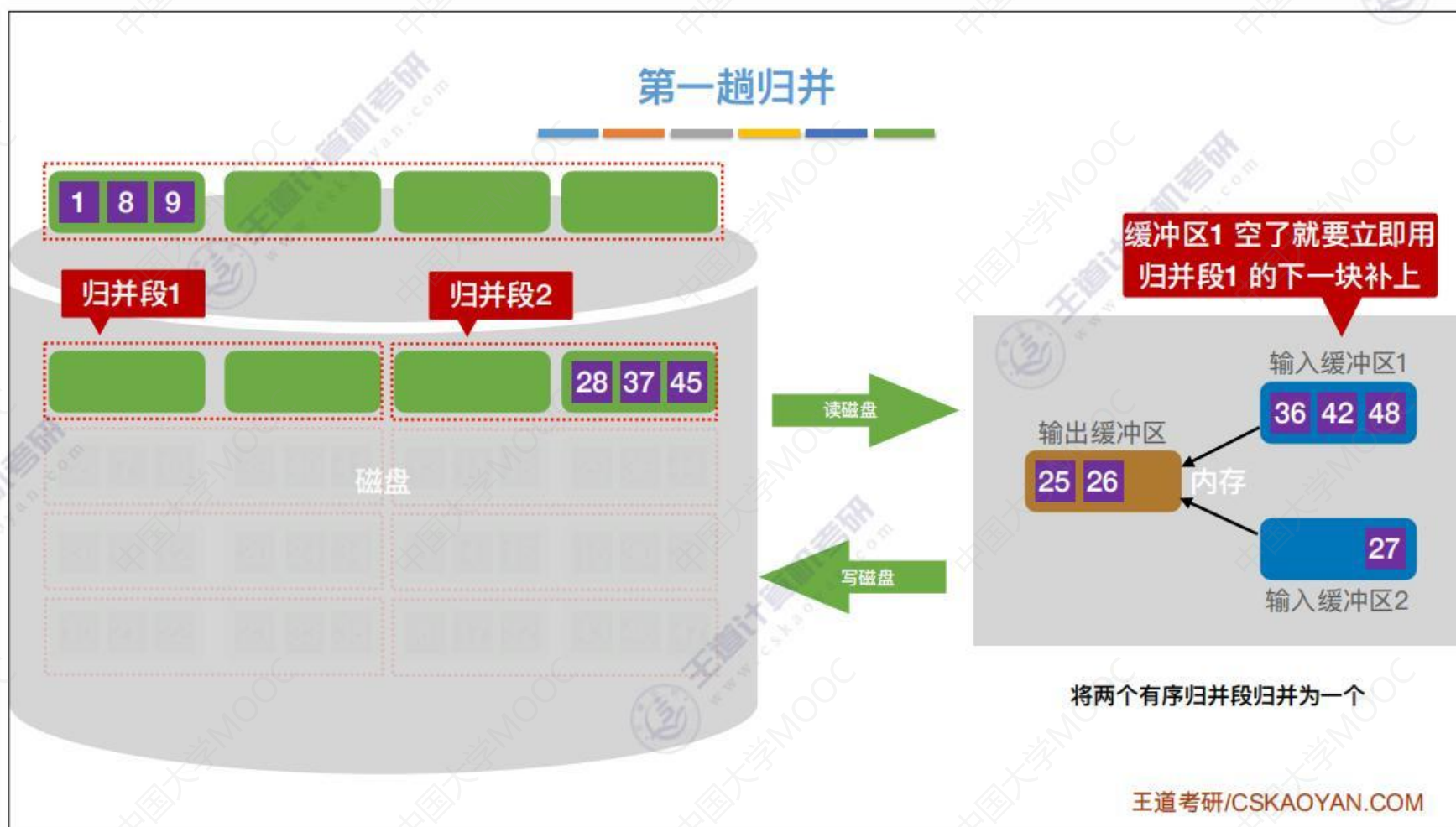
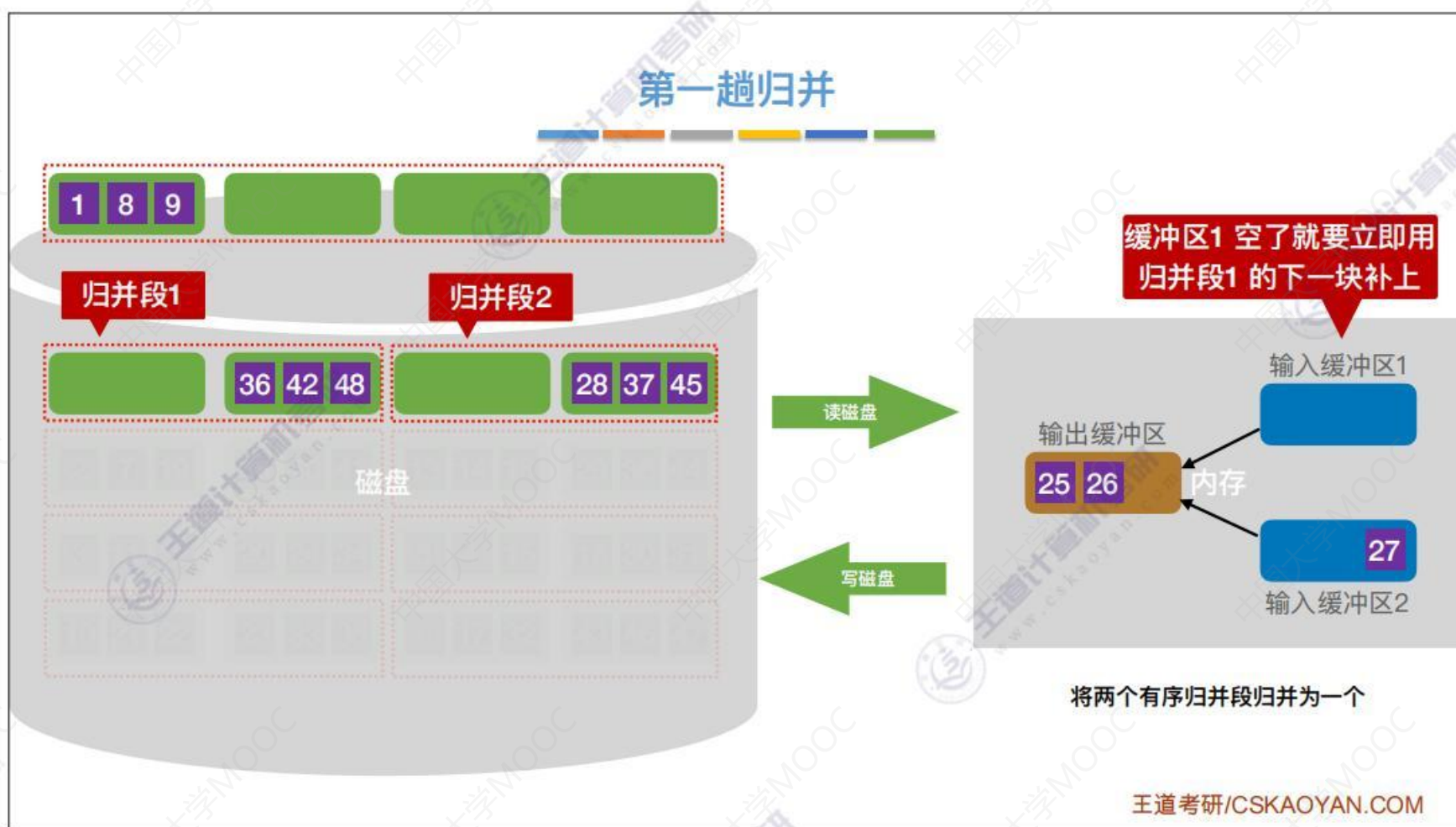
把8个有序子序列（初始归并段）两两归并

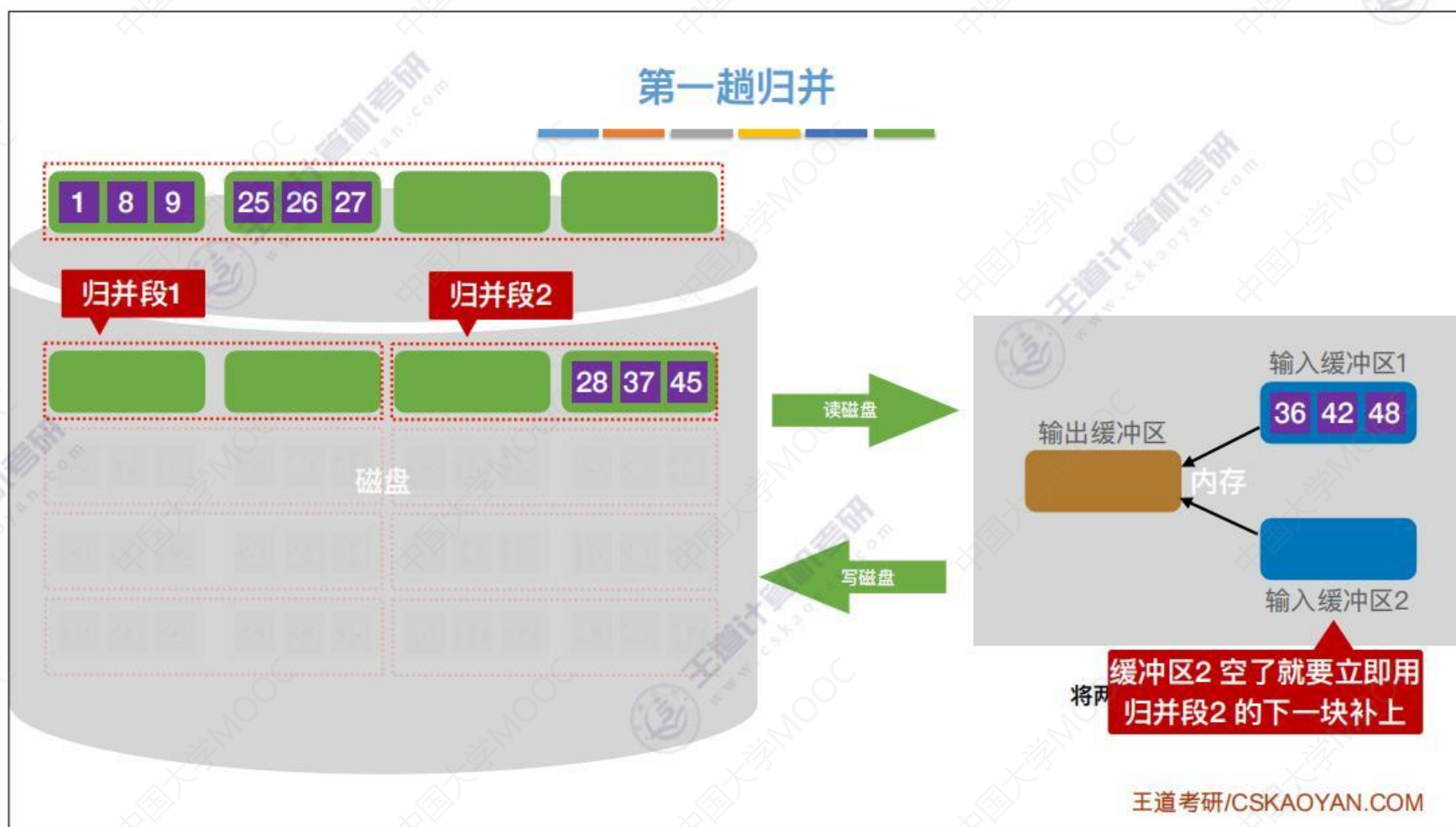
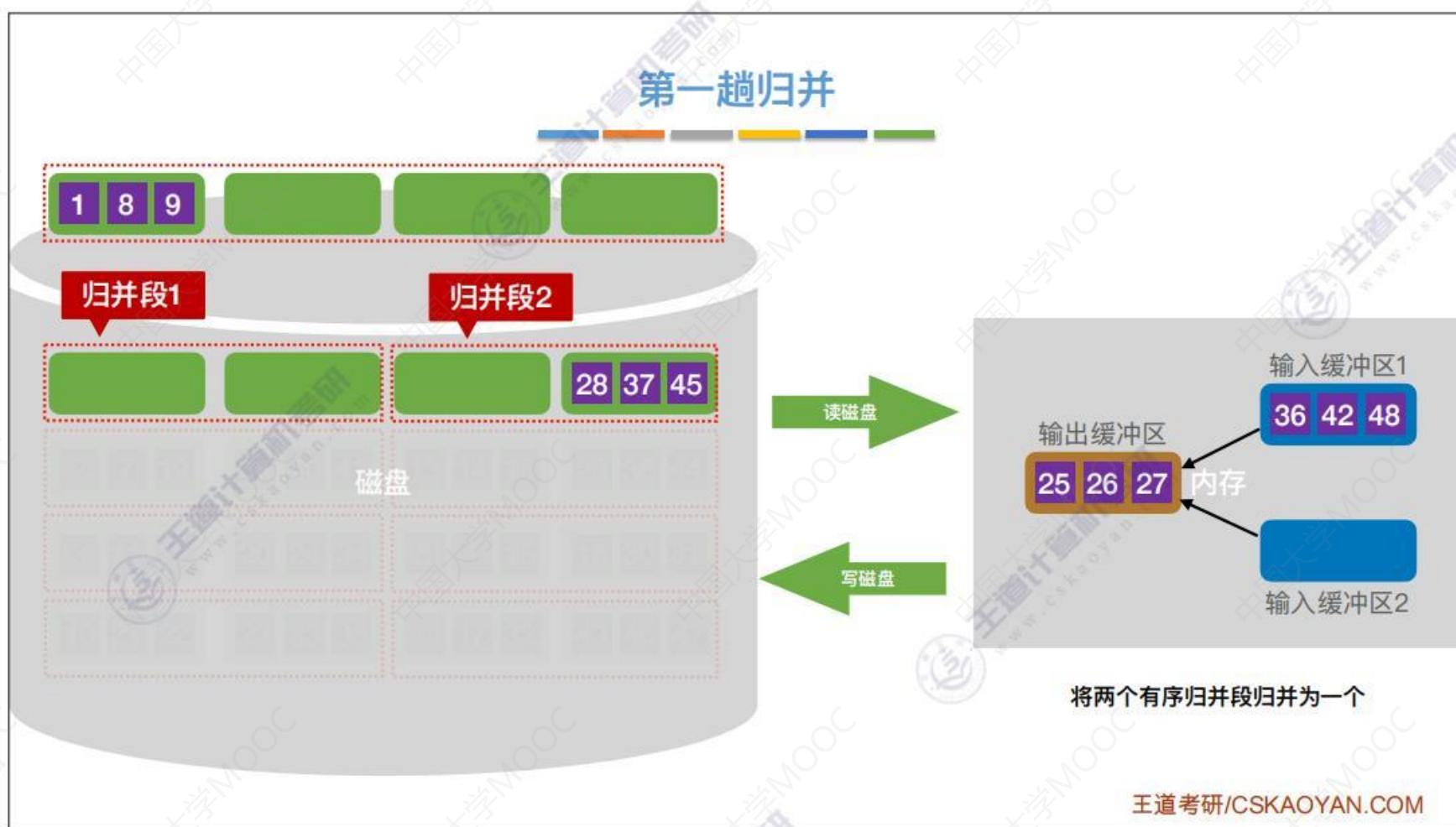


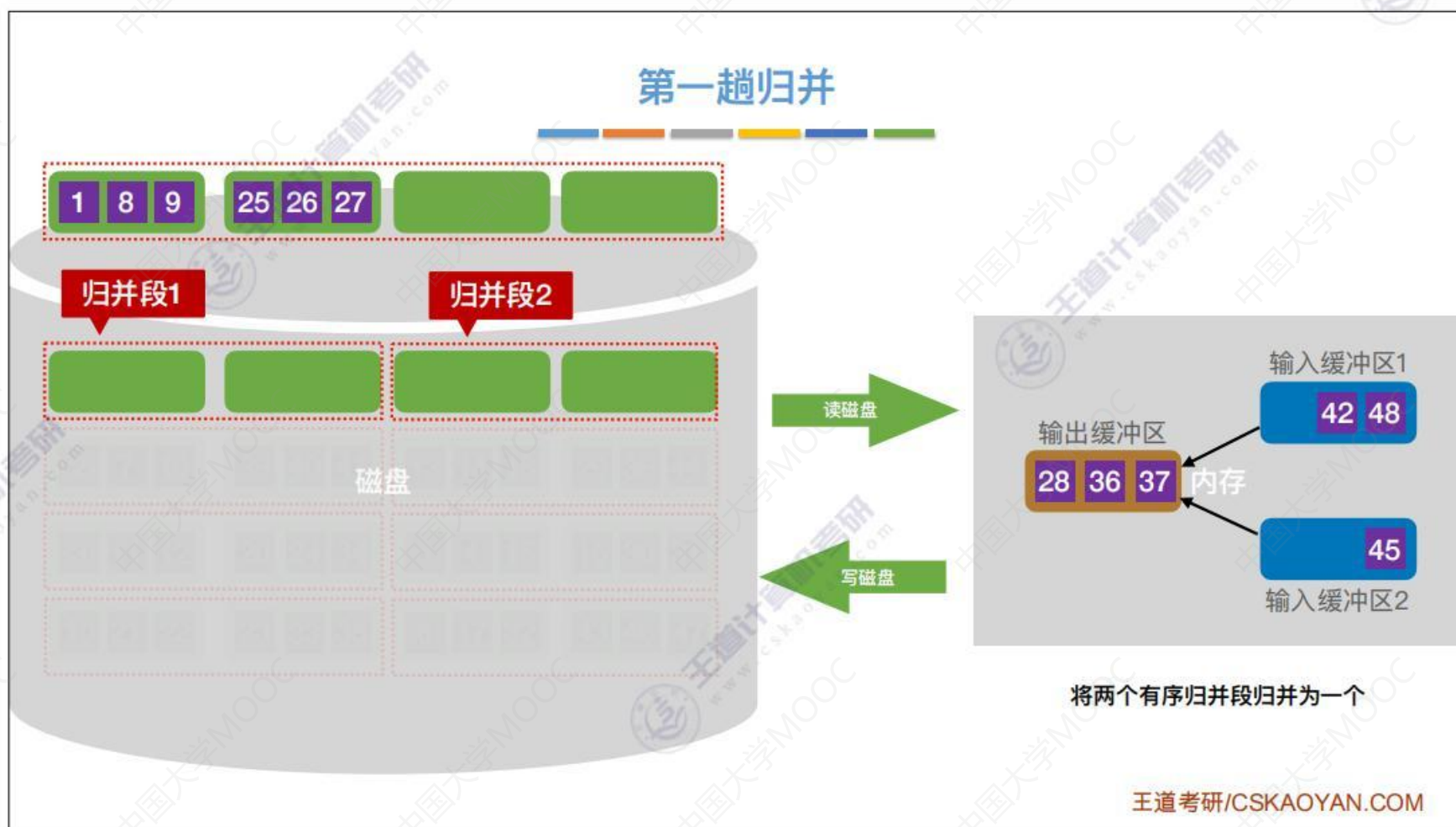
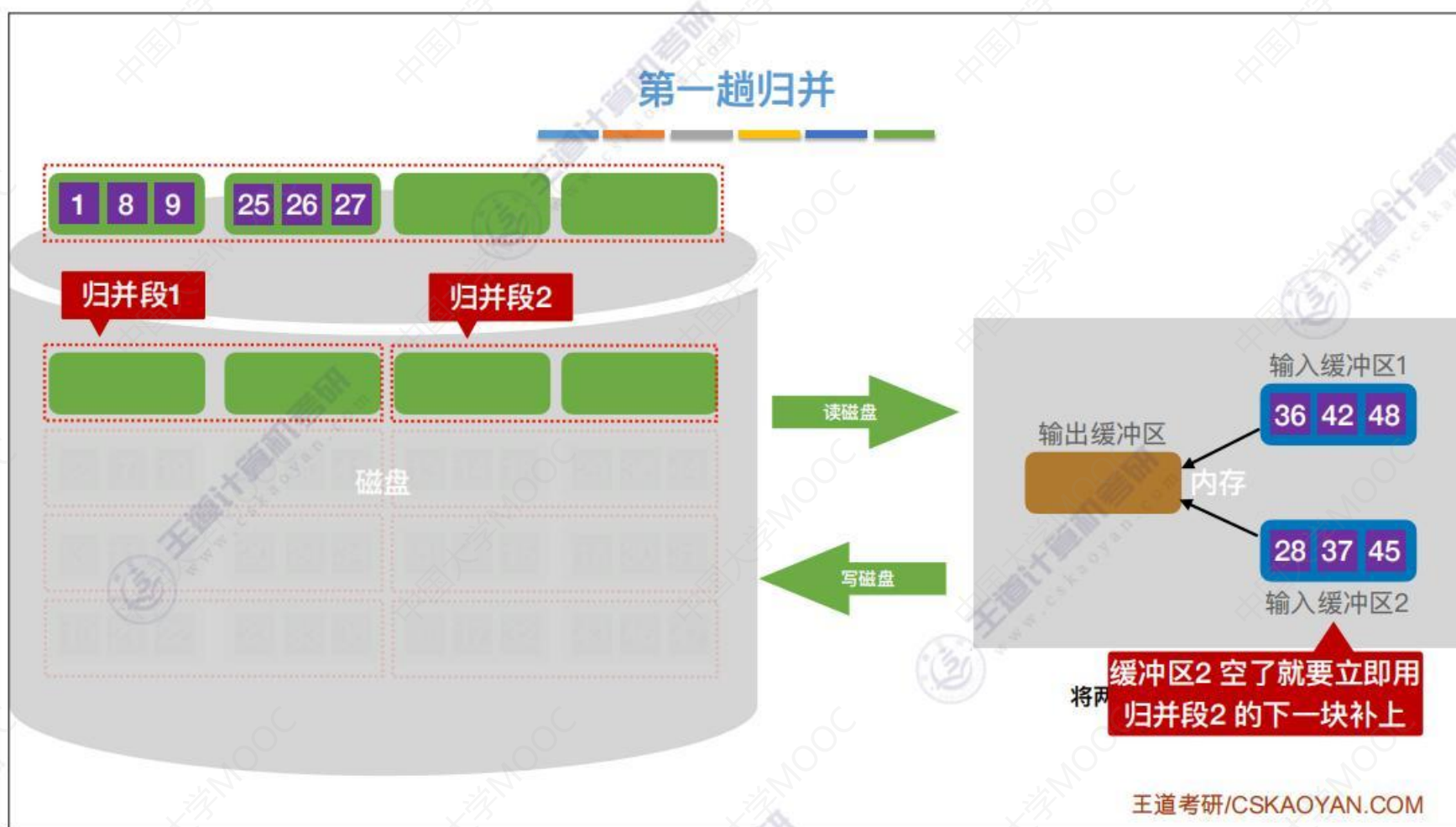


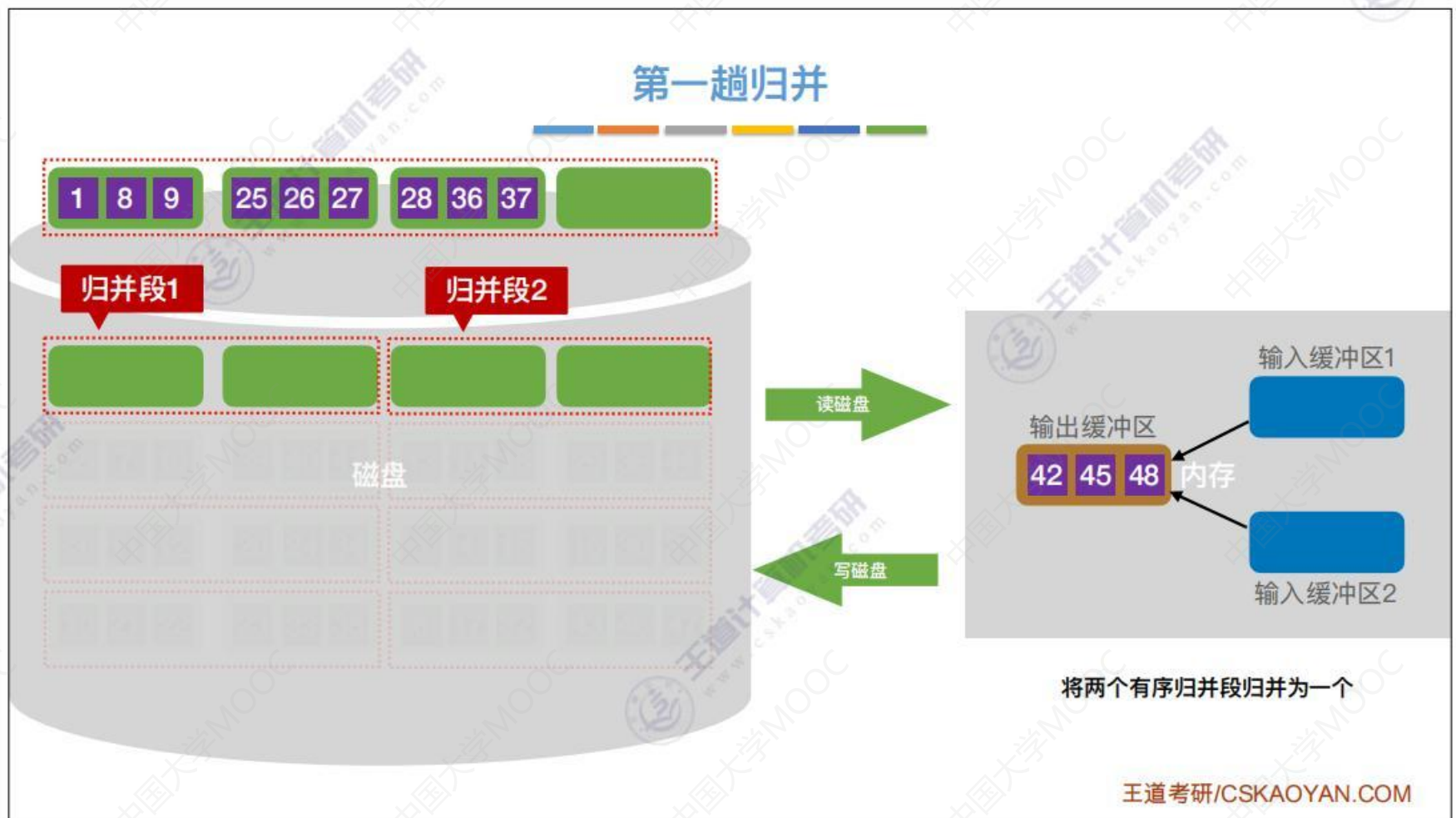
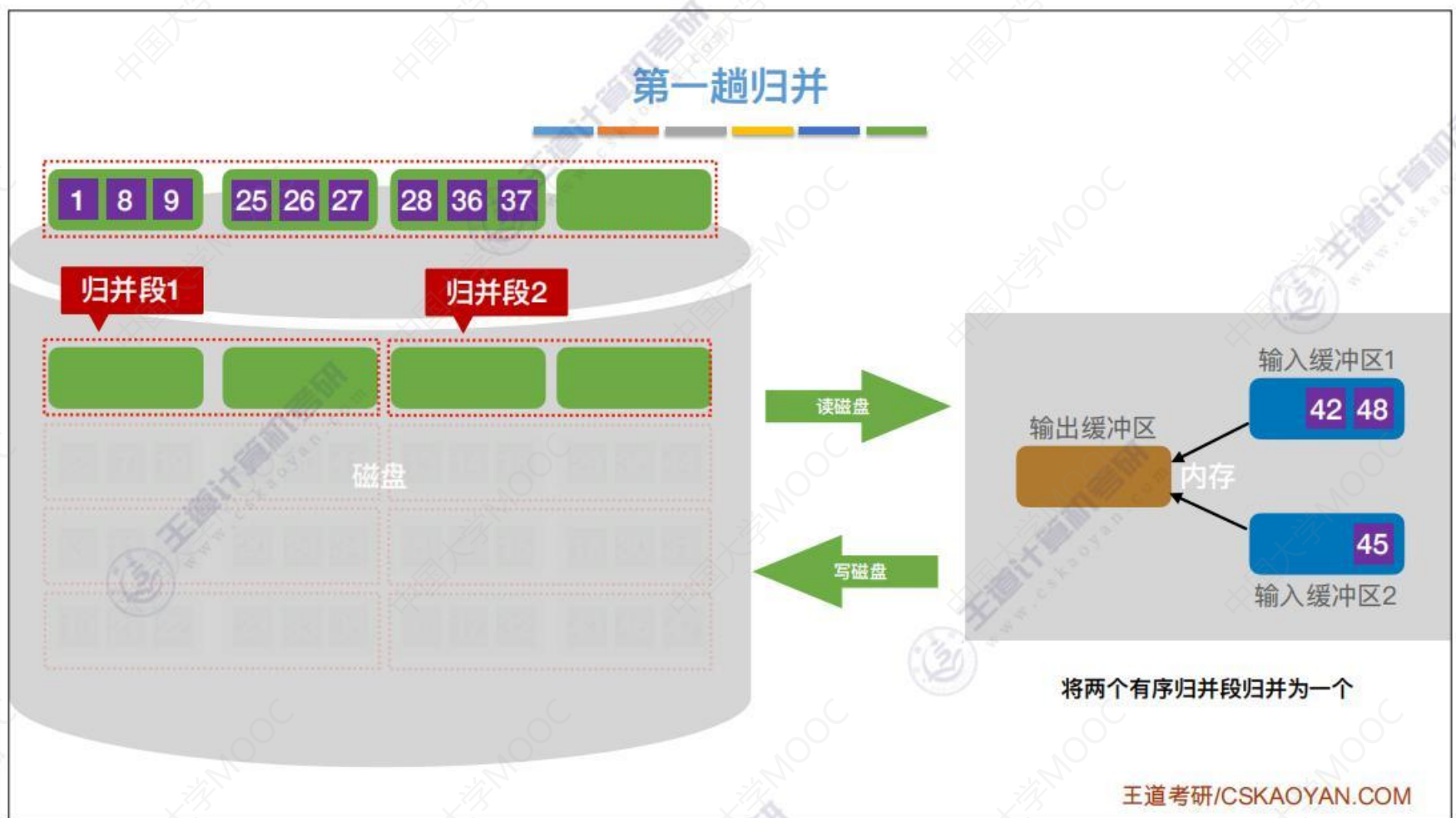


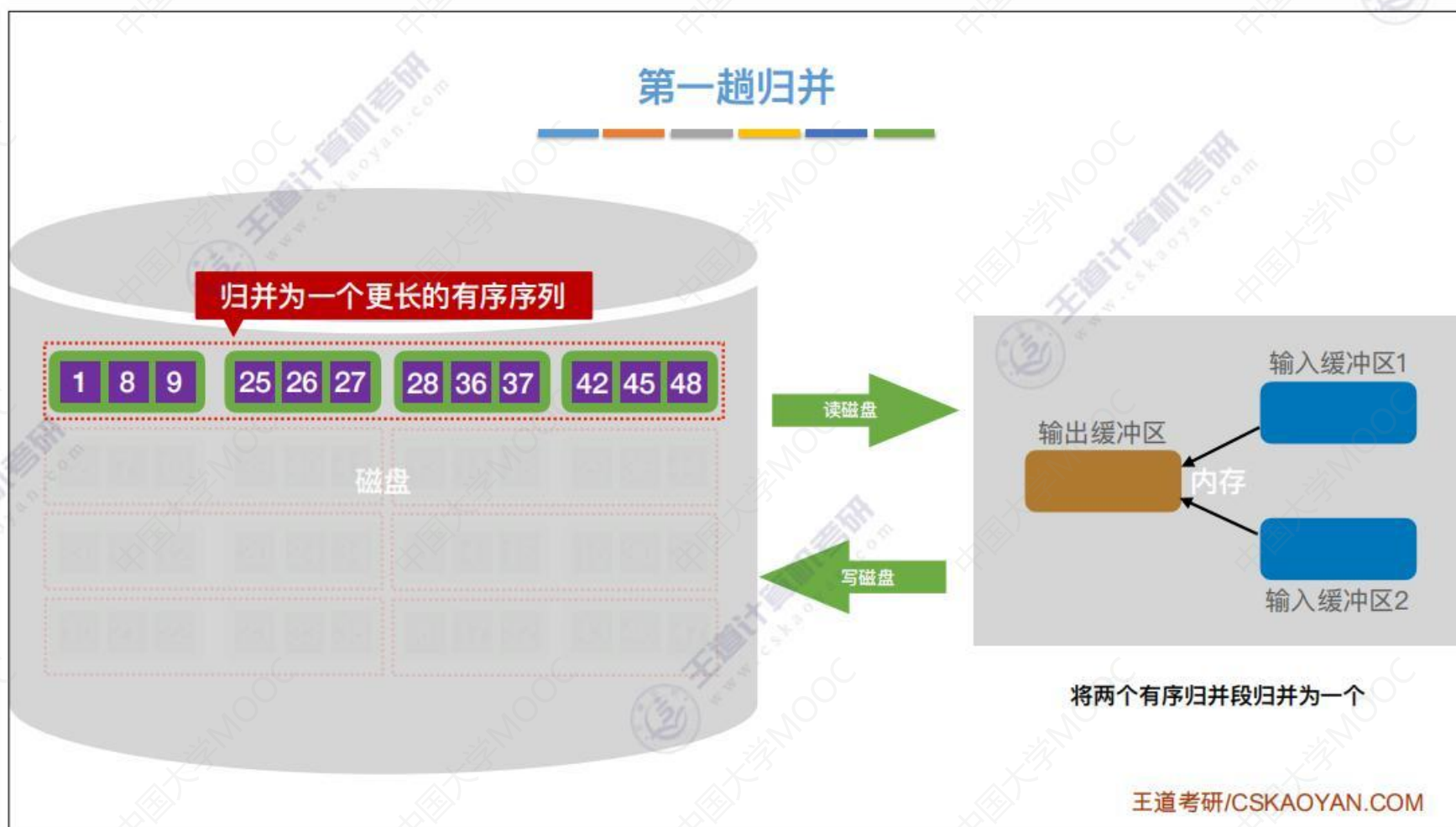
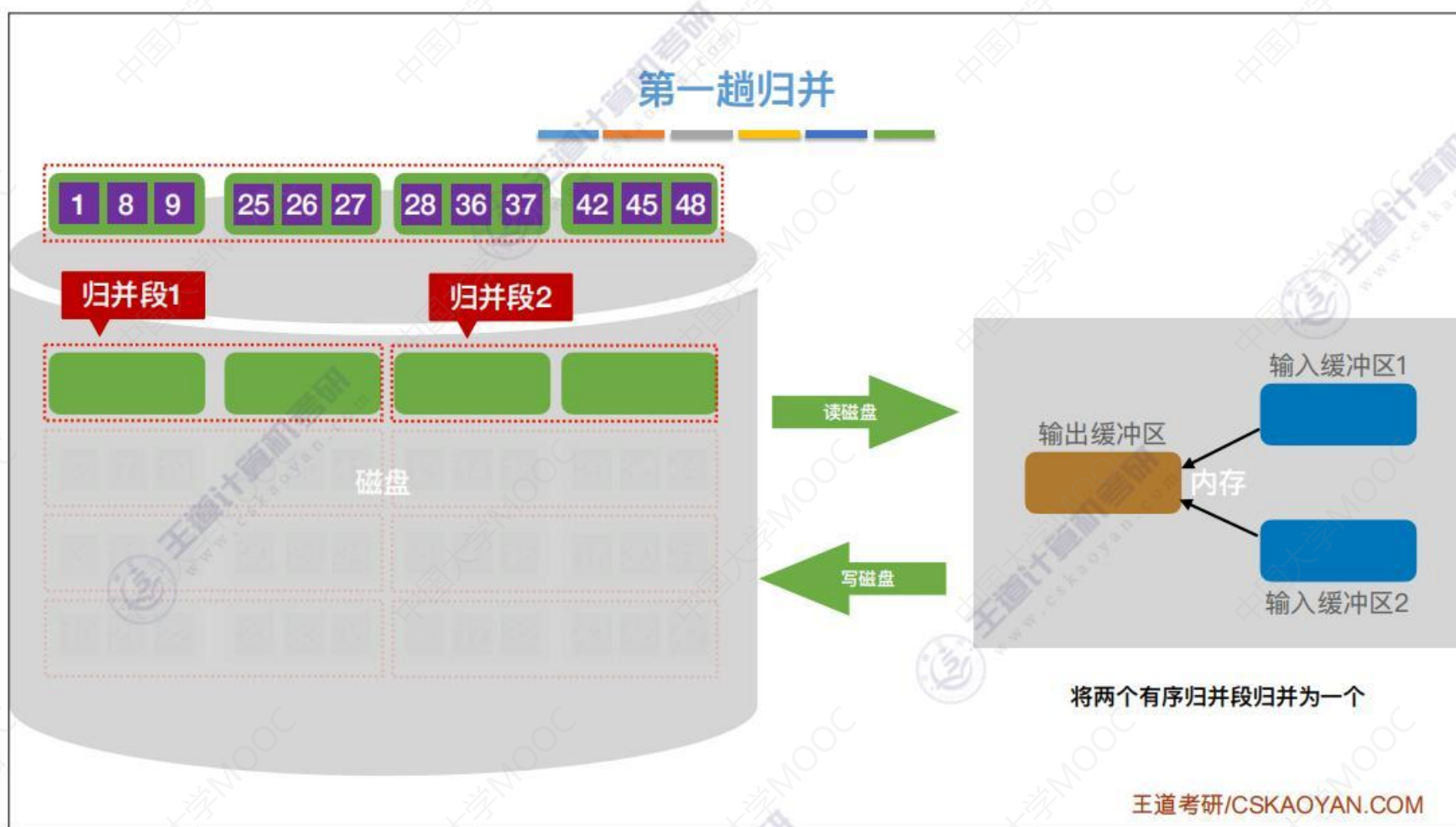


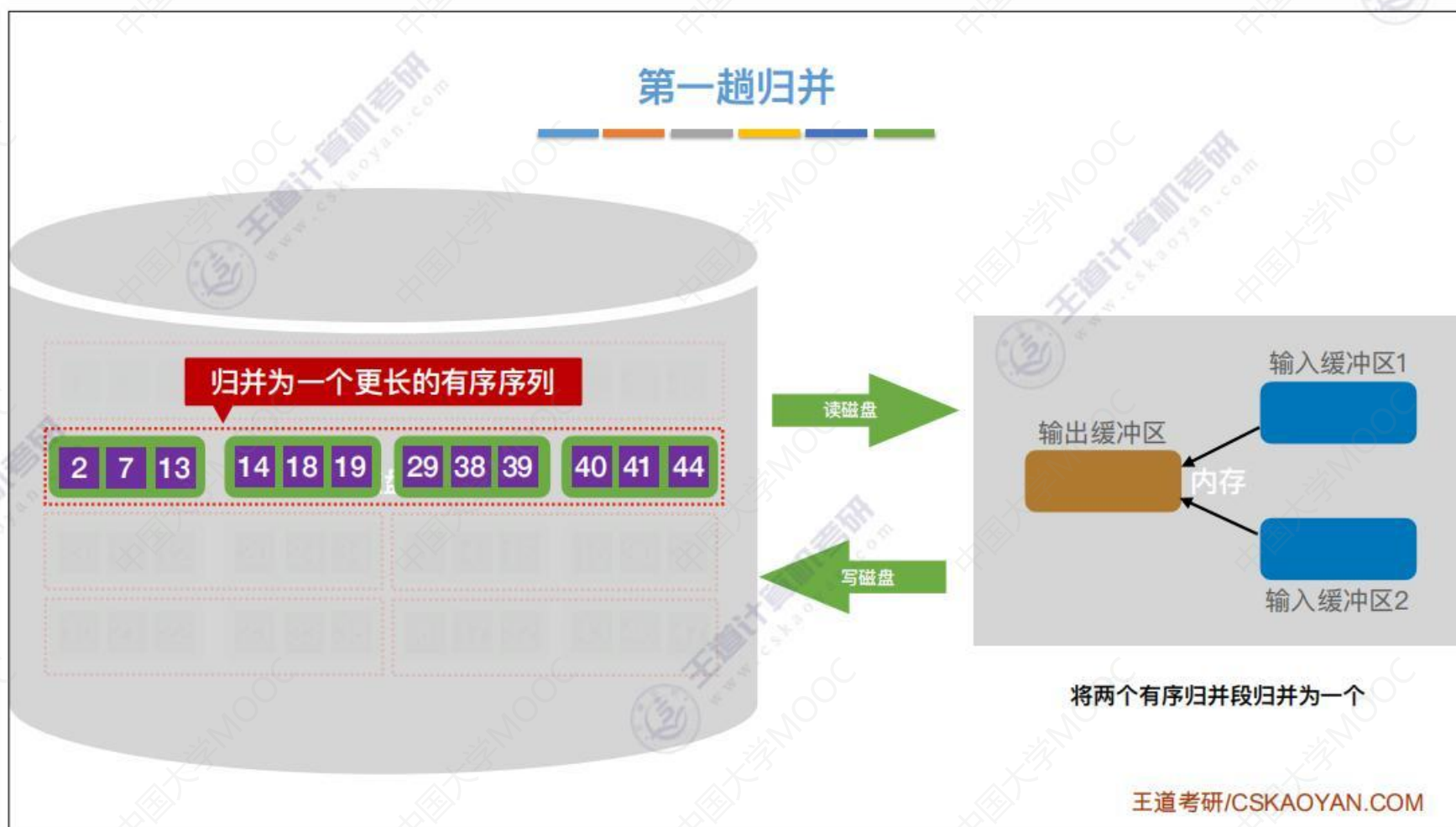
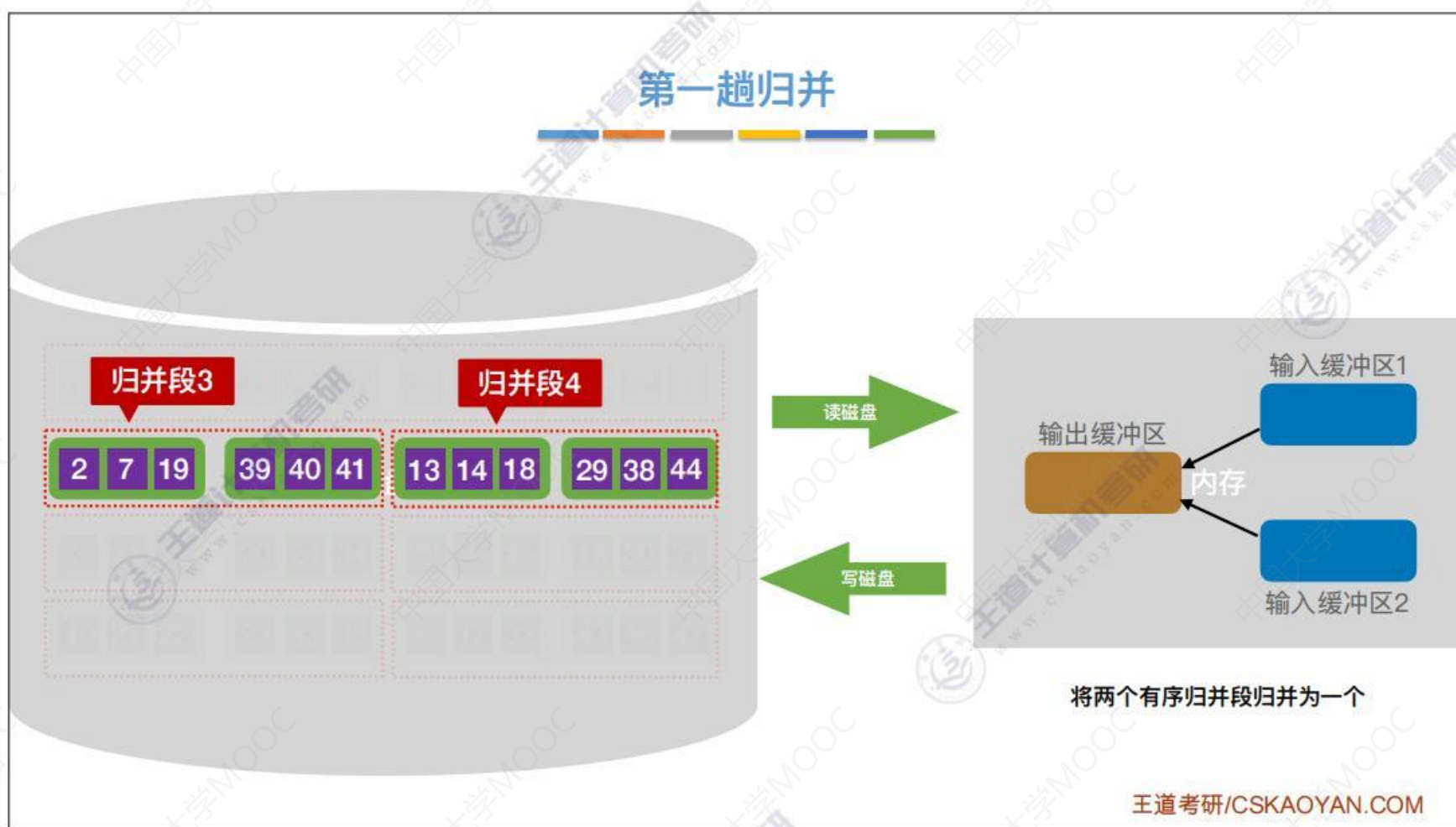


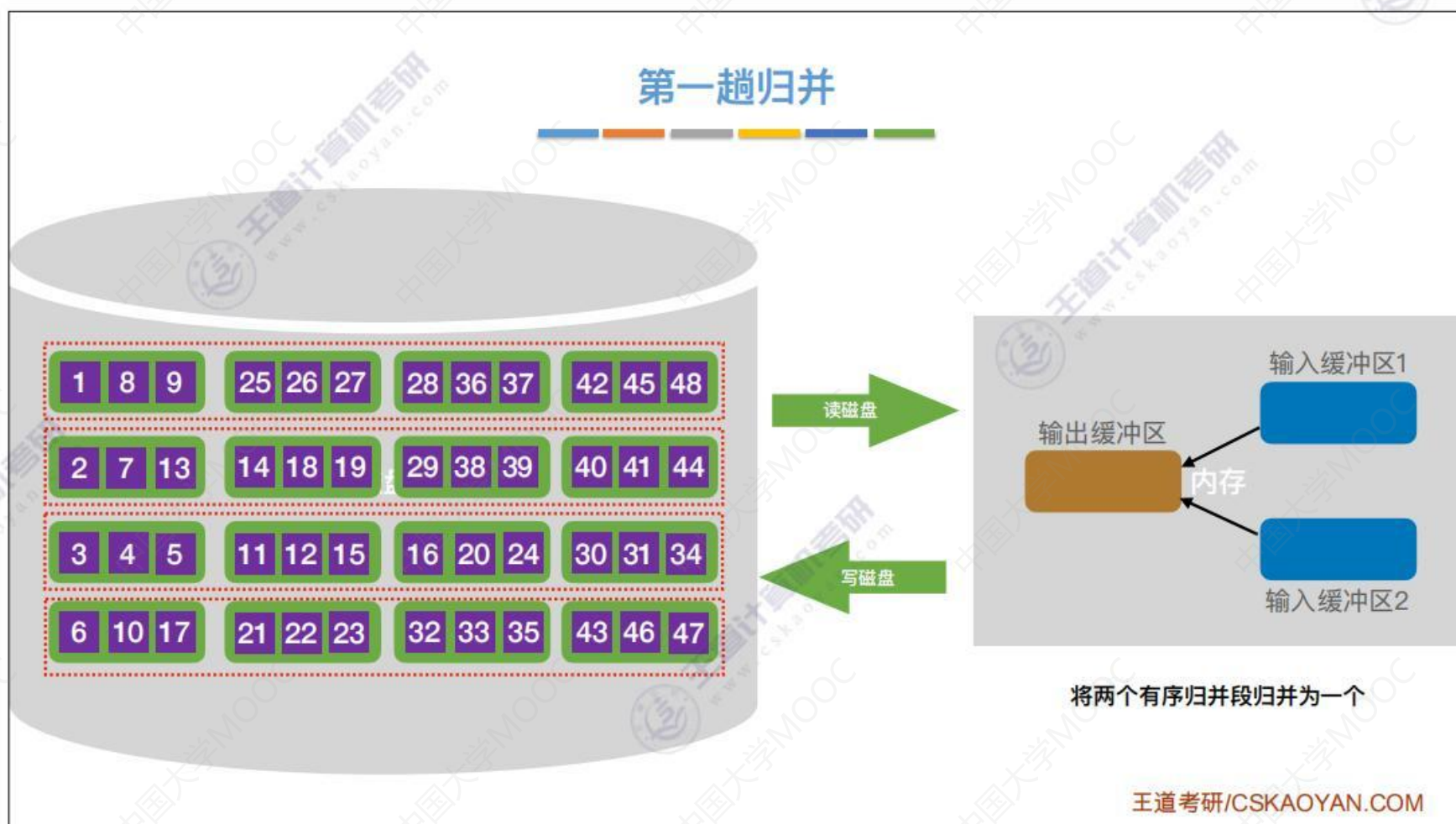
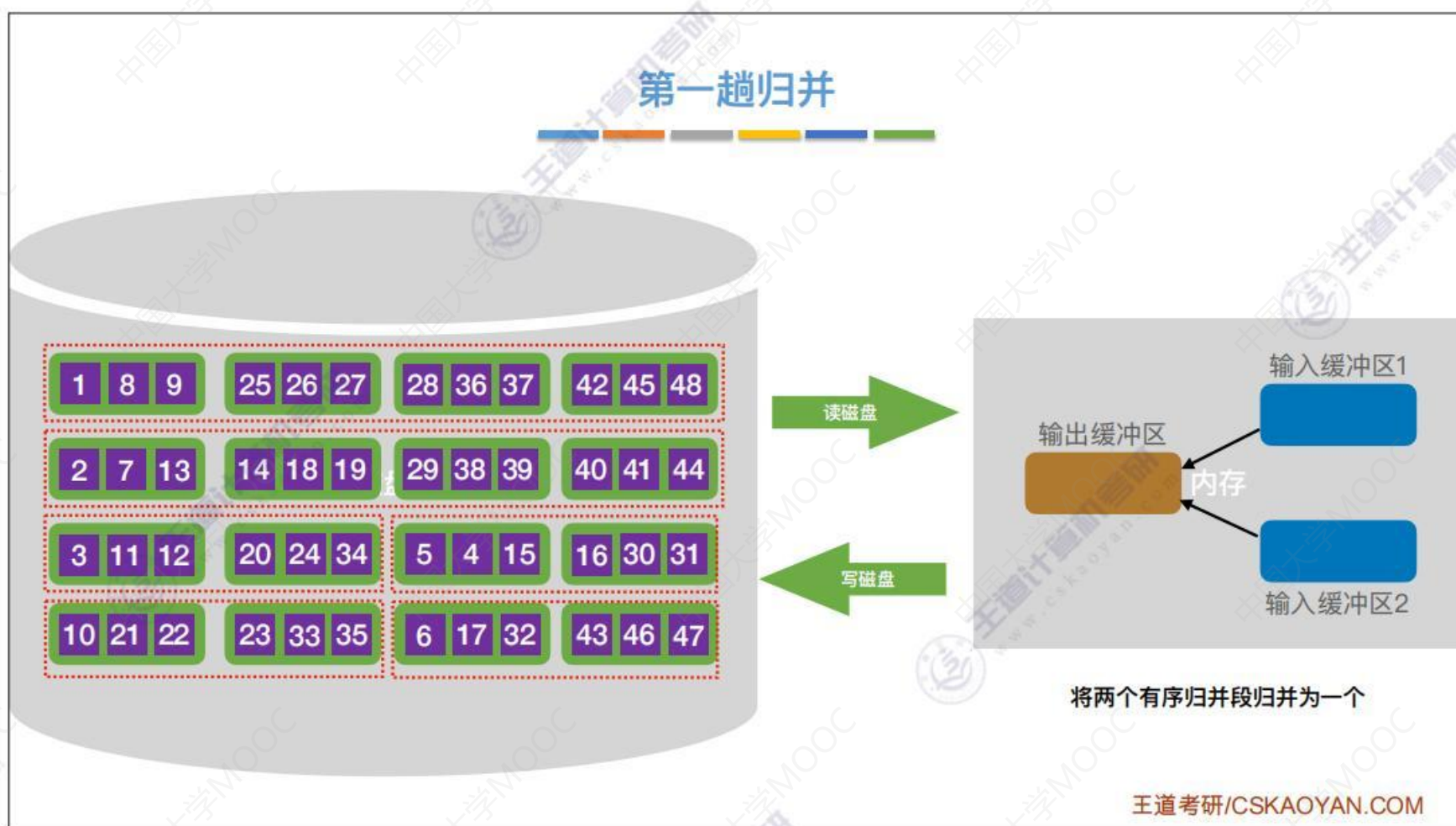


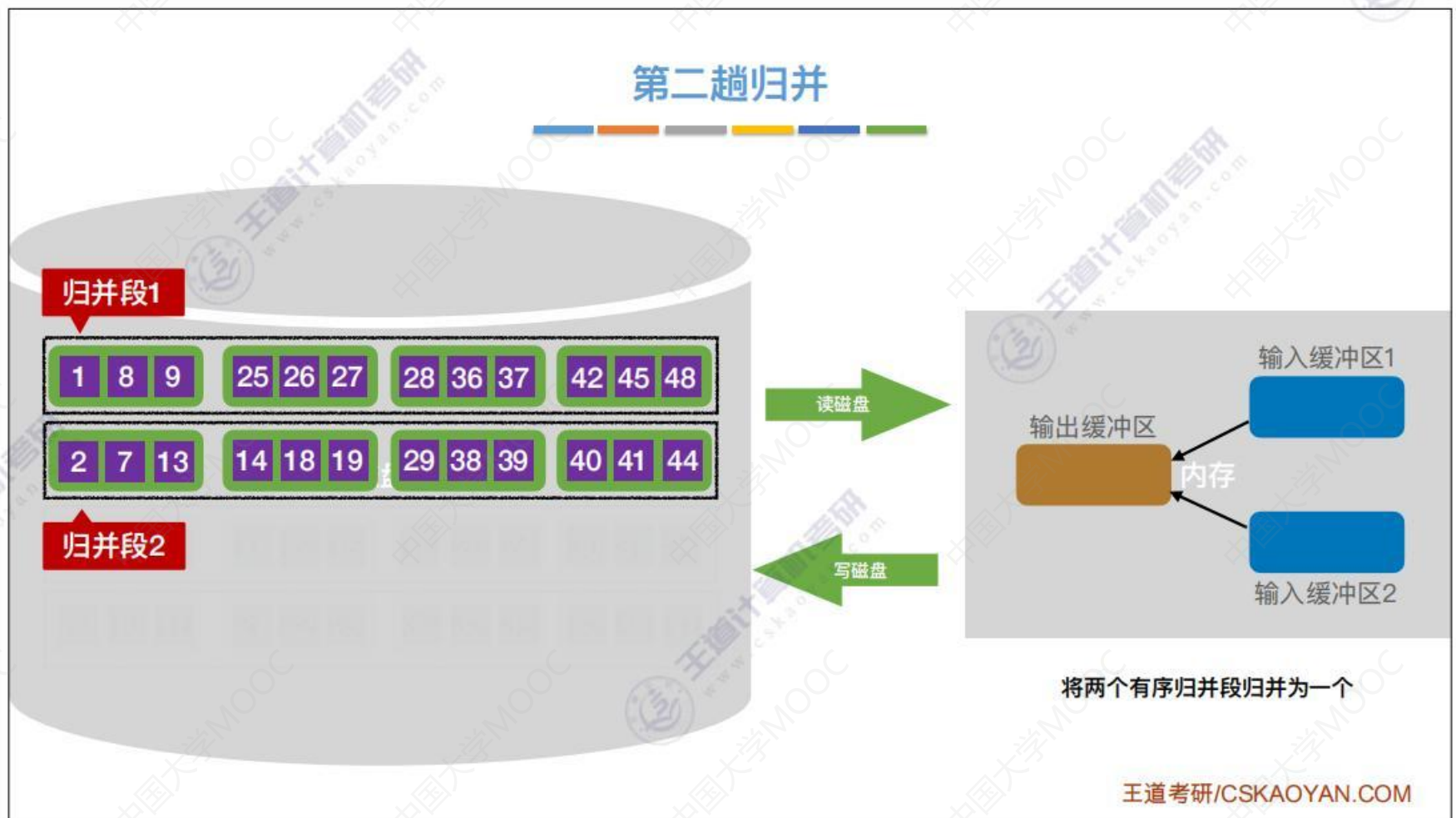
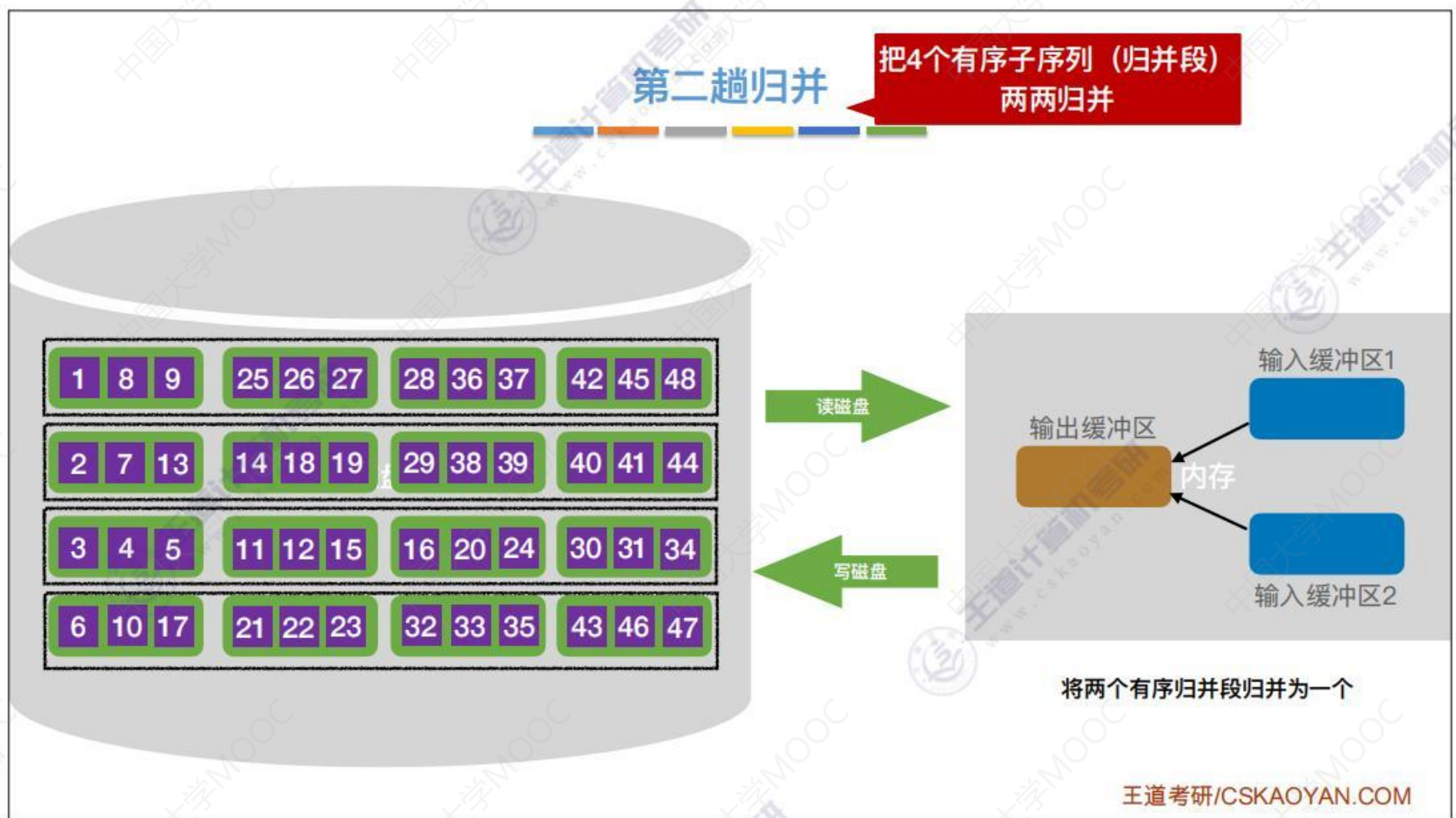


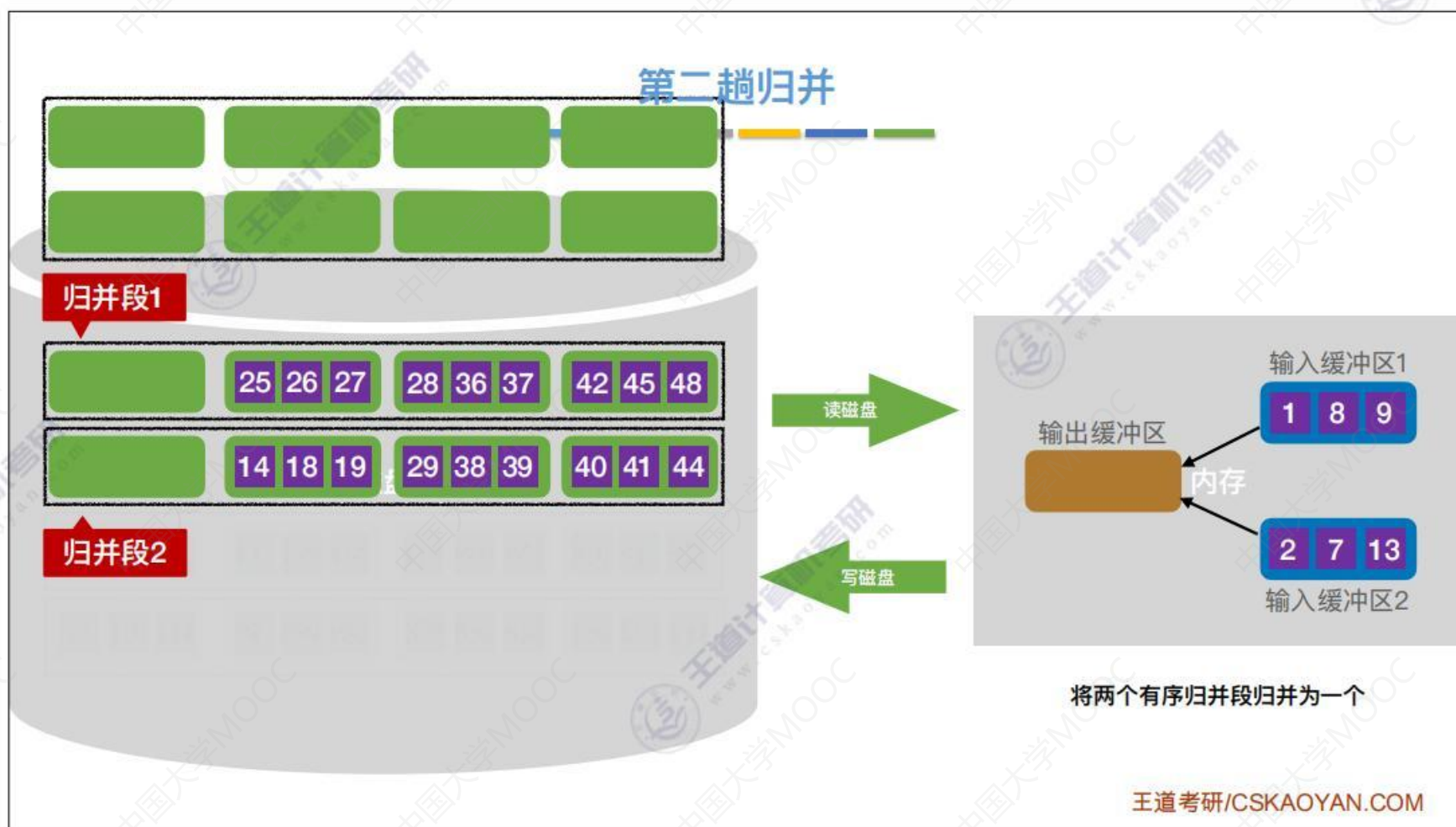
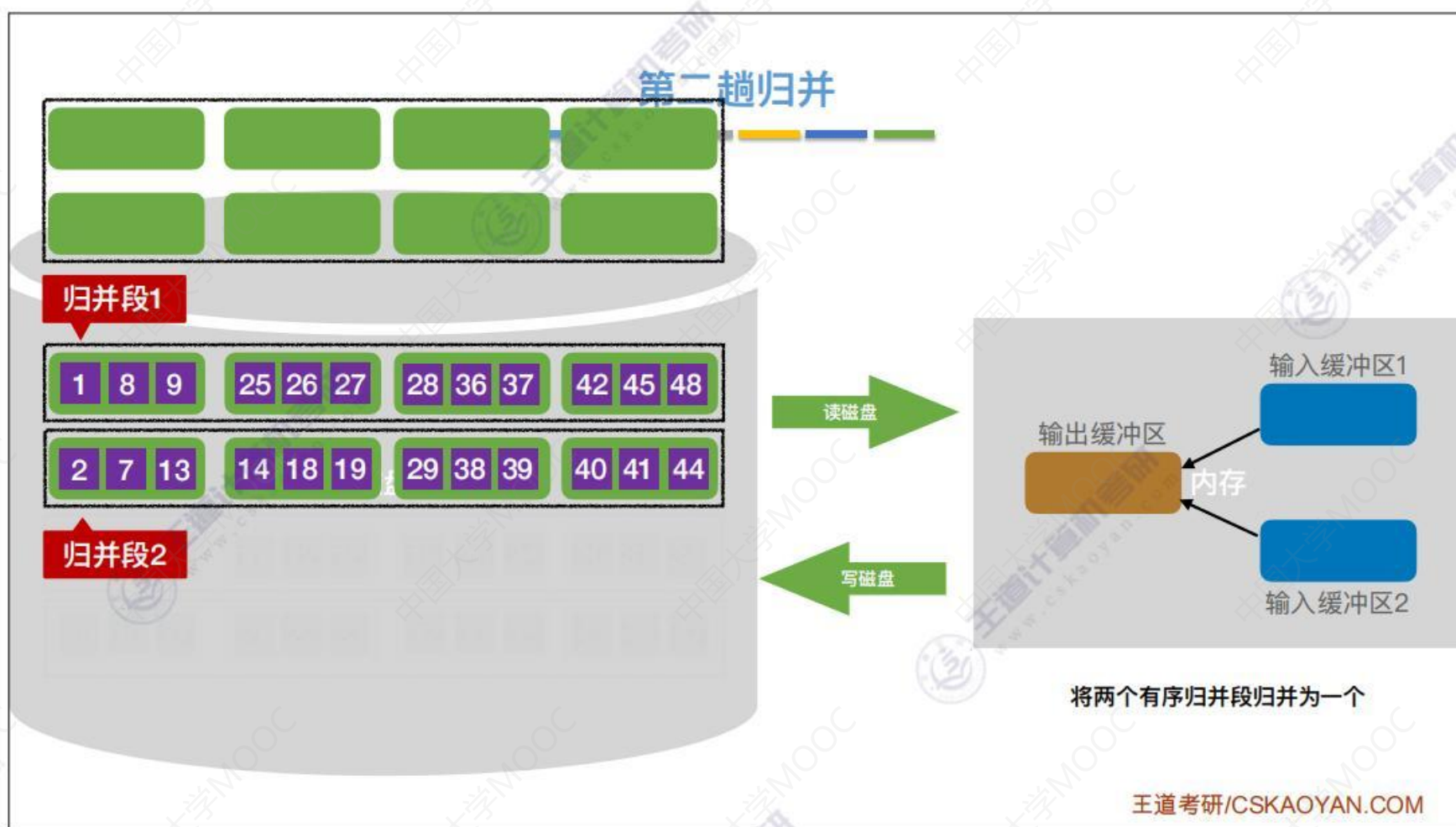


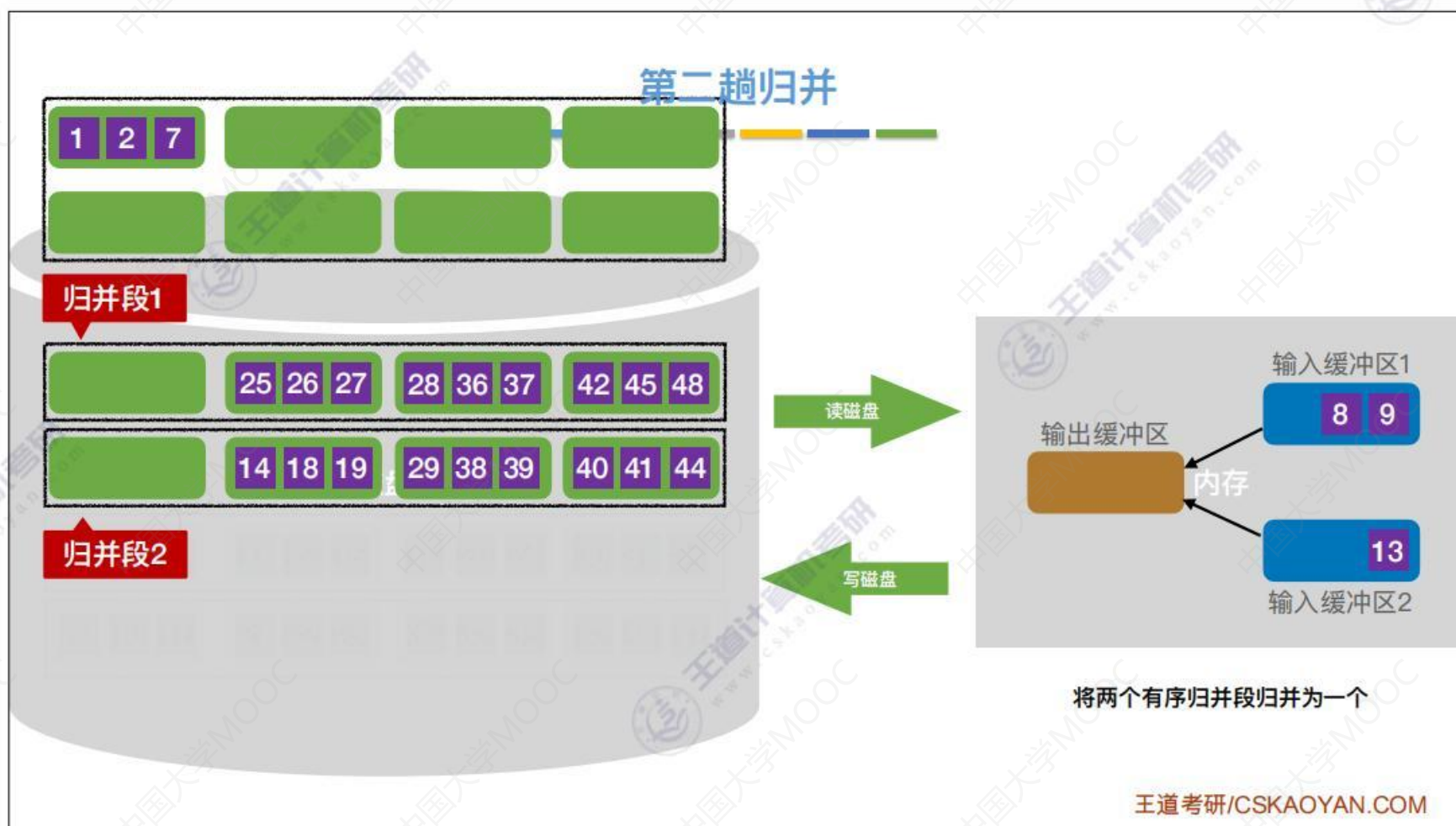
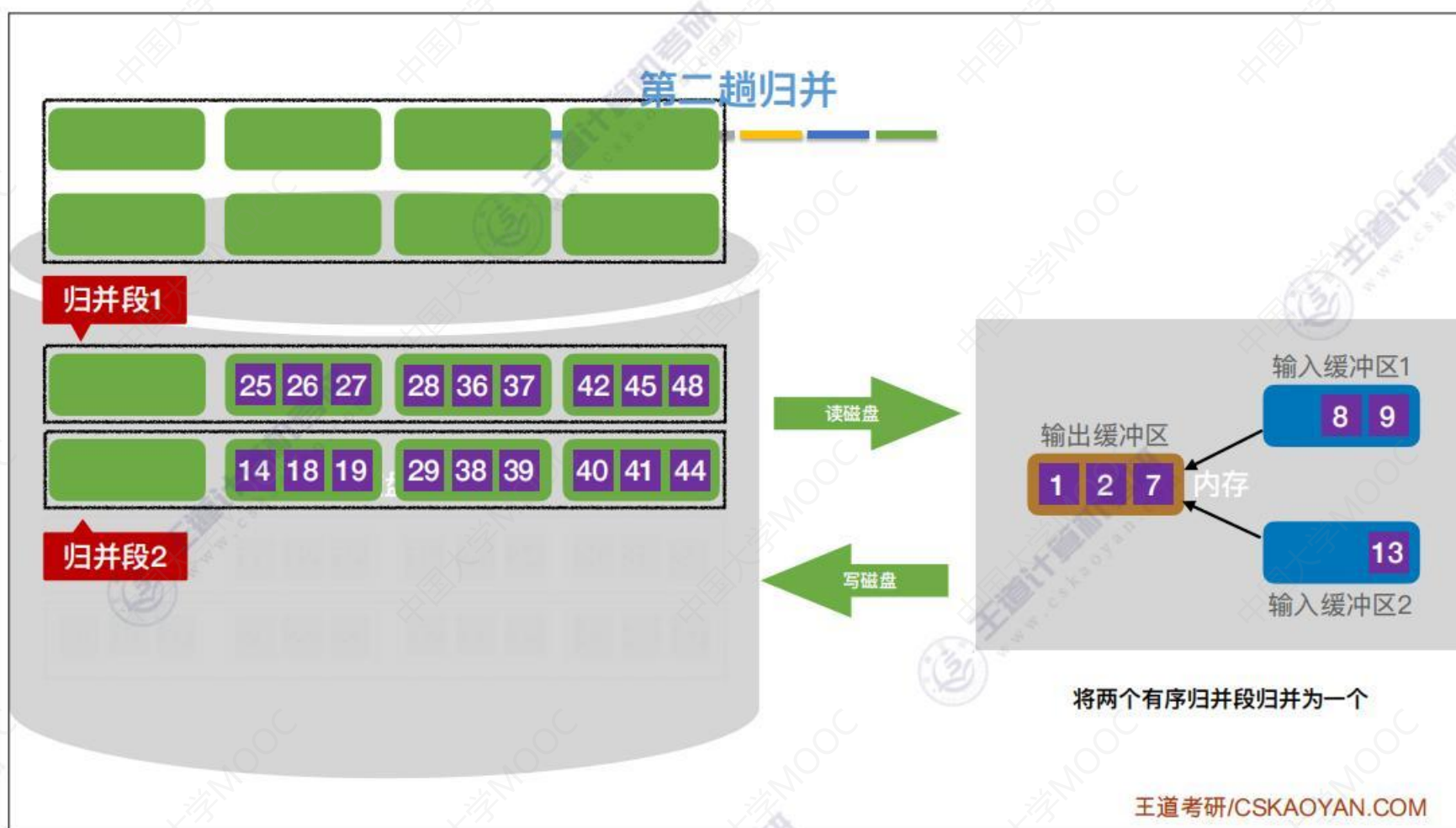


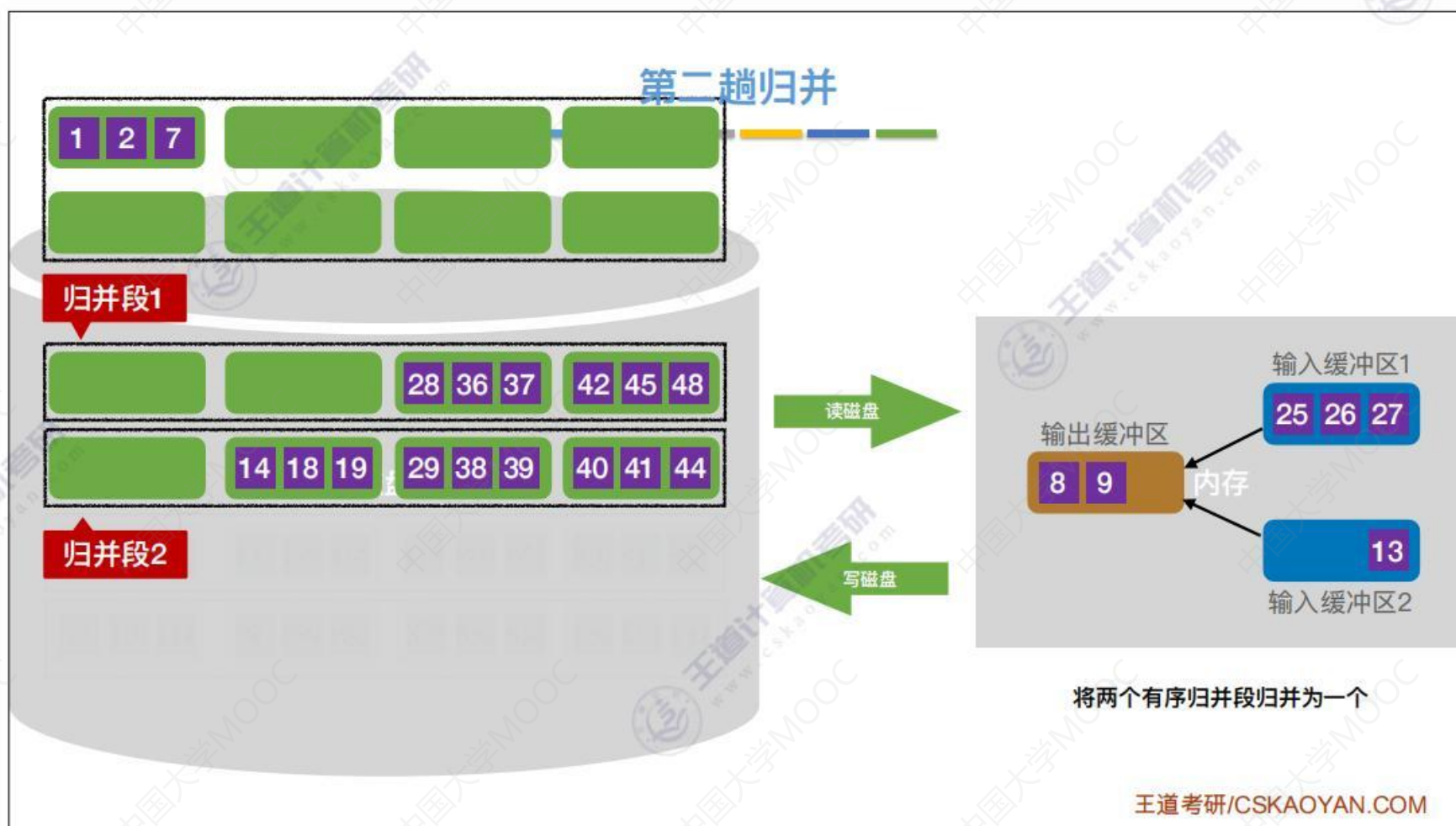
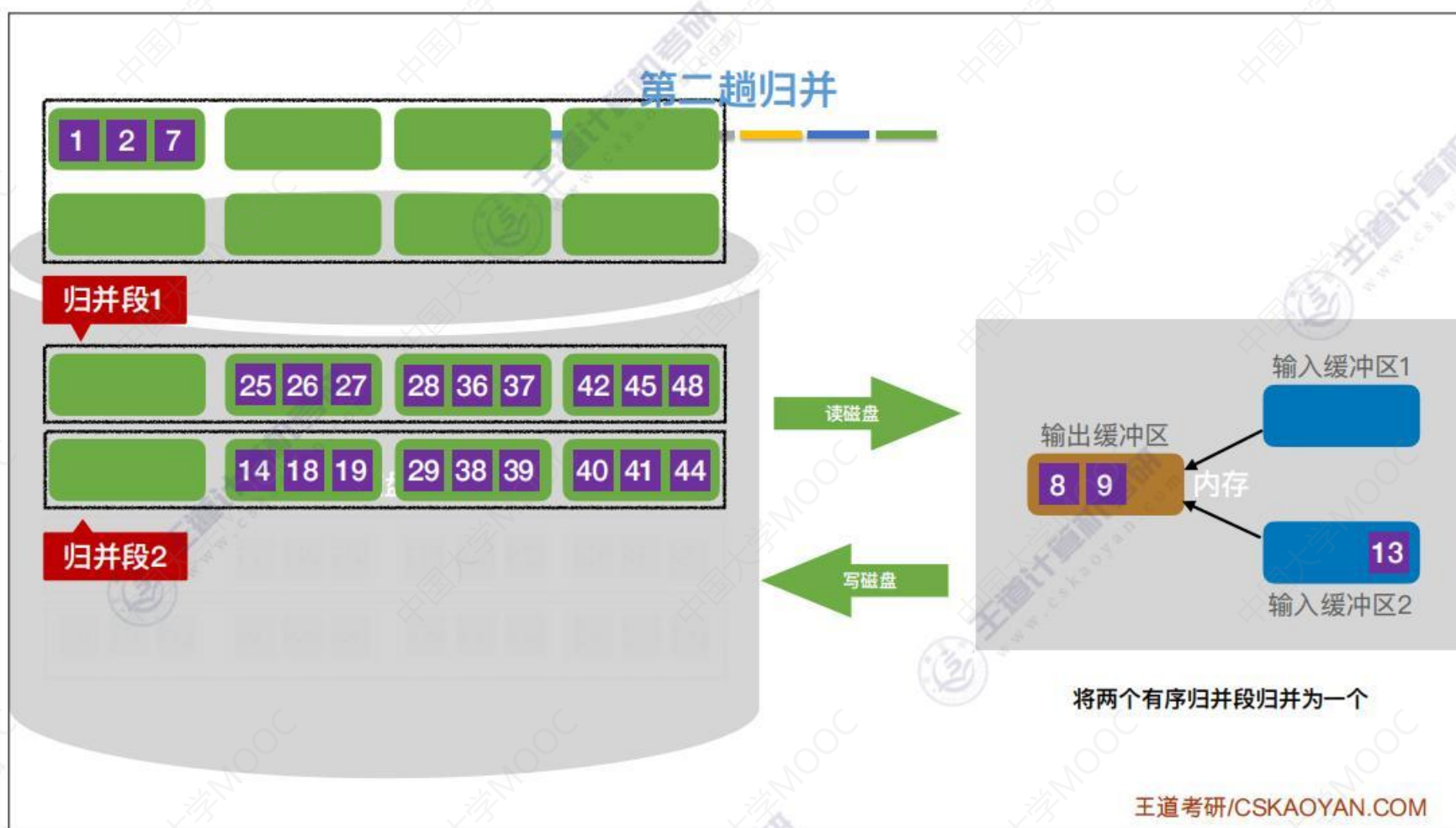


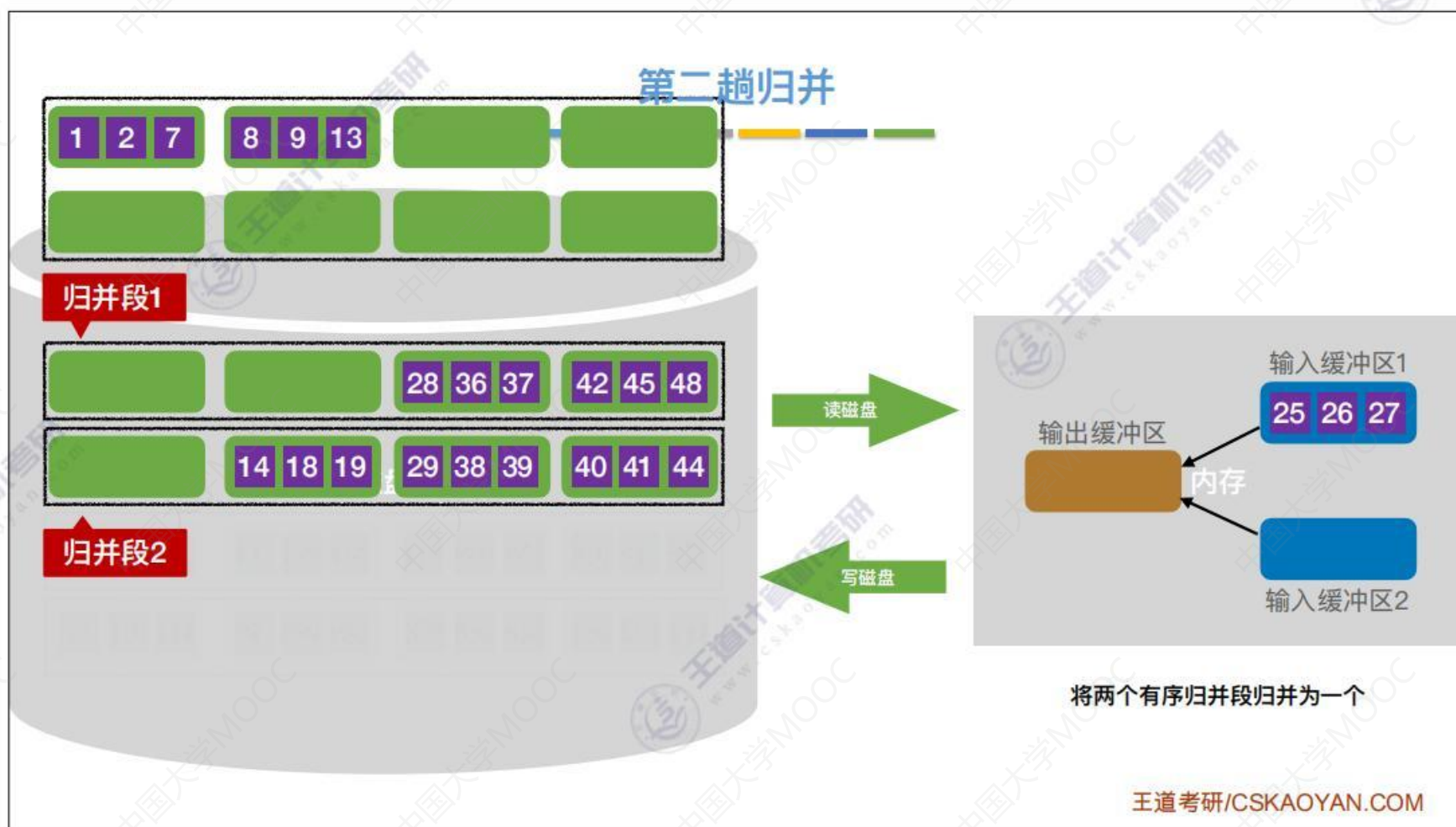
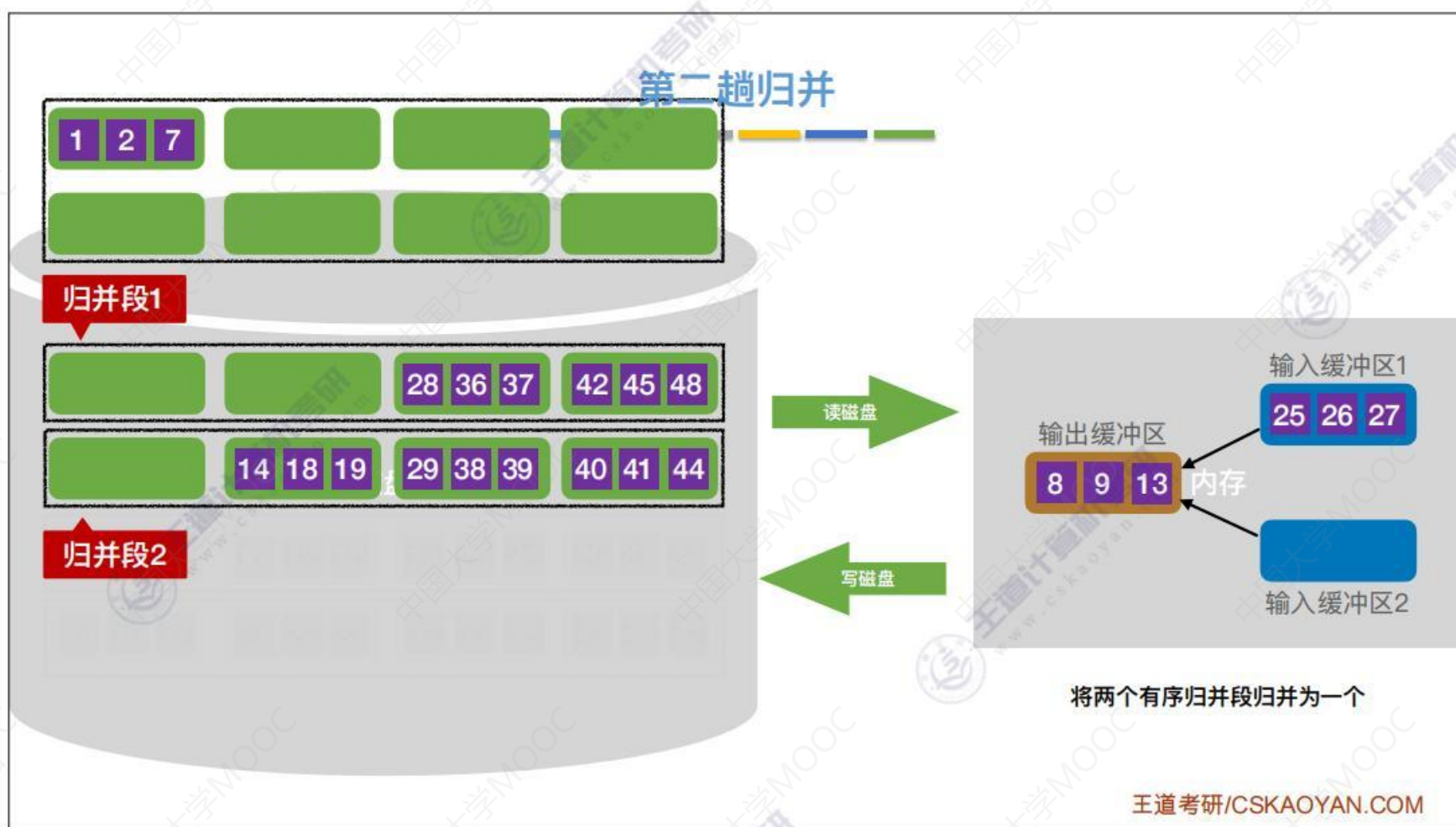


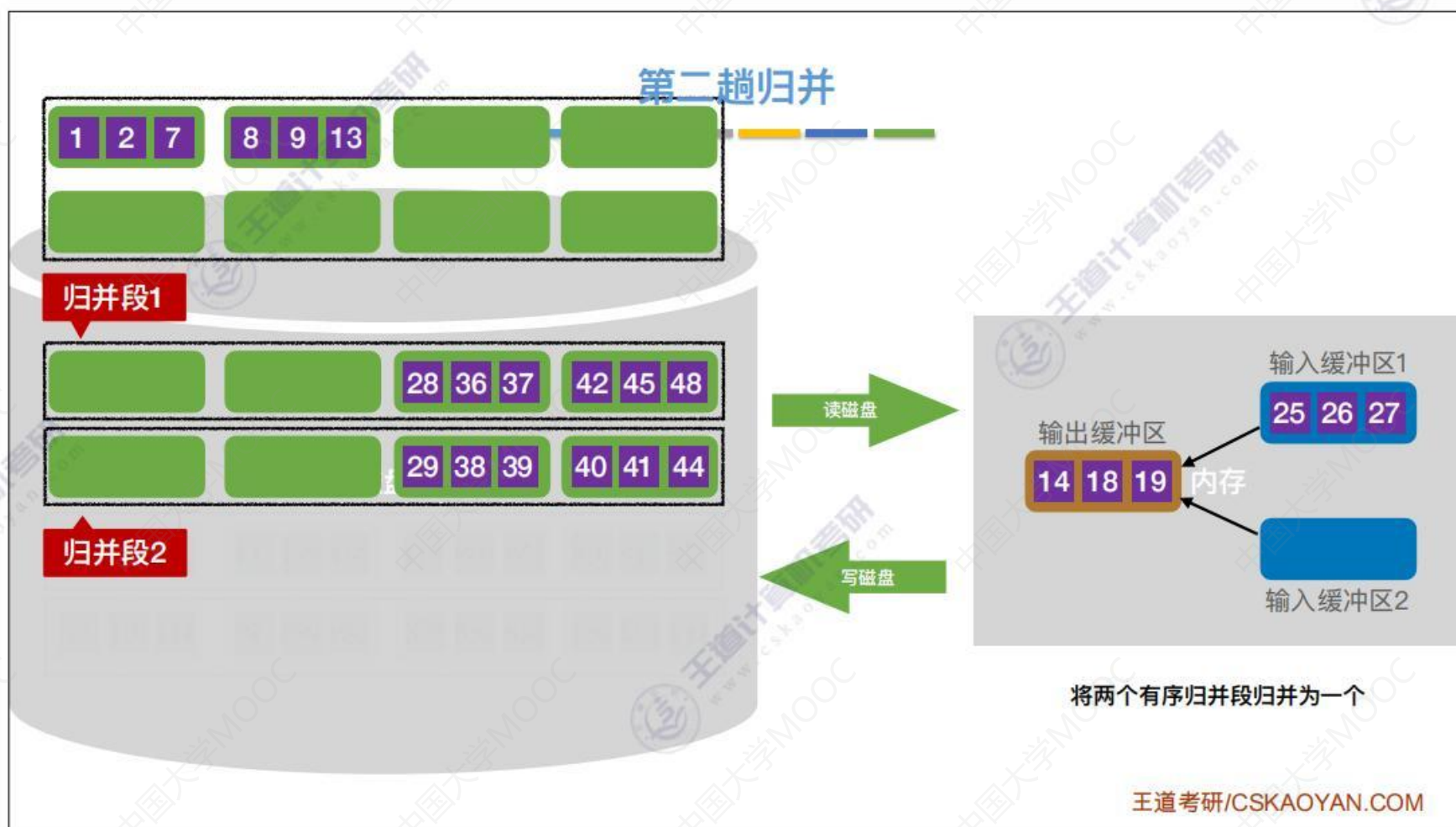
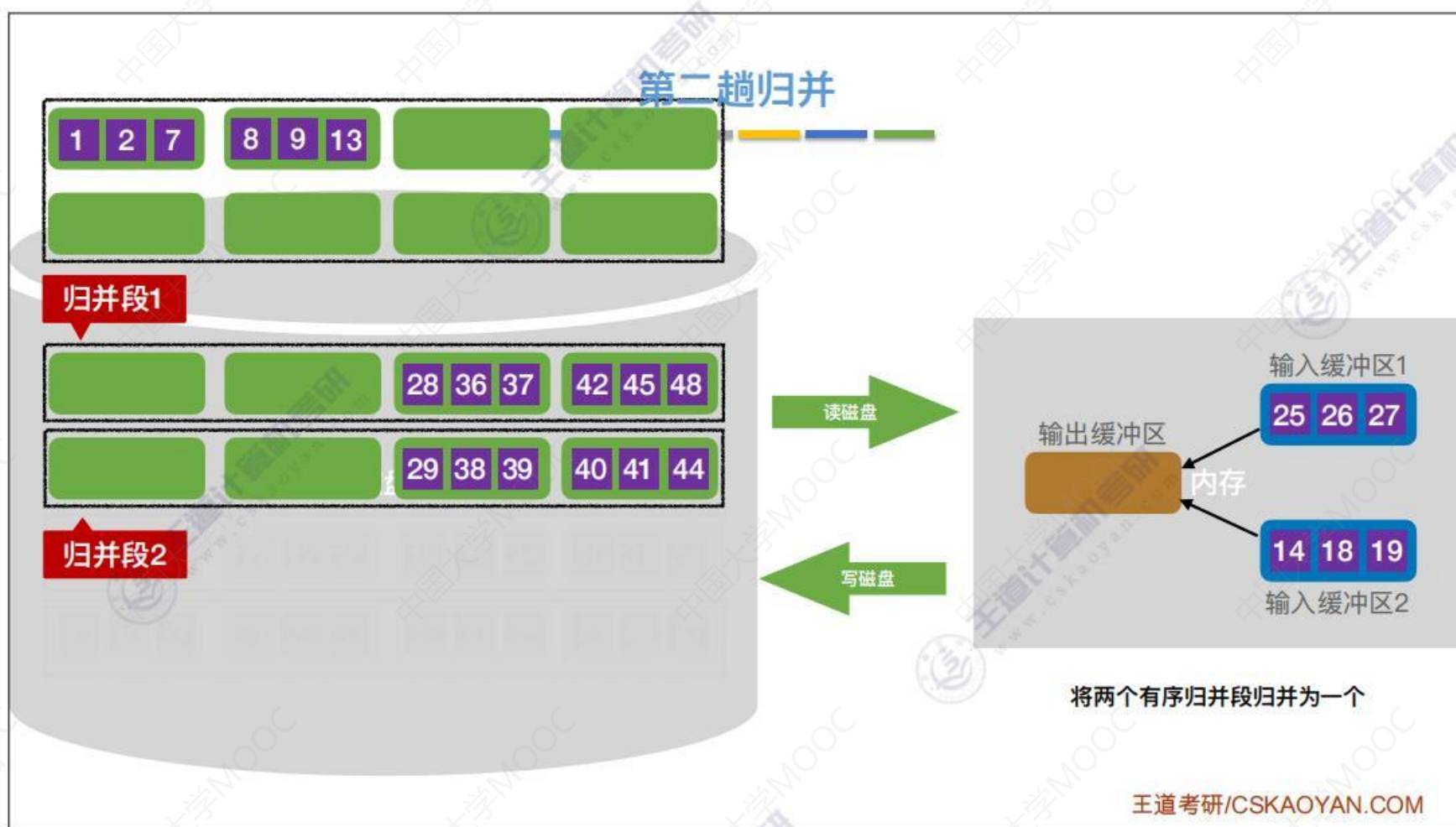


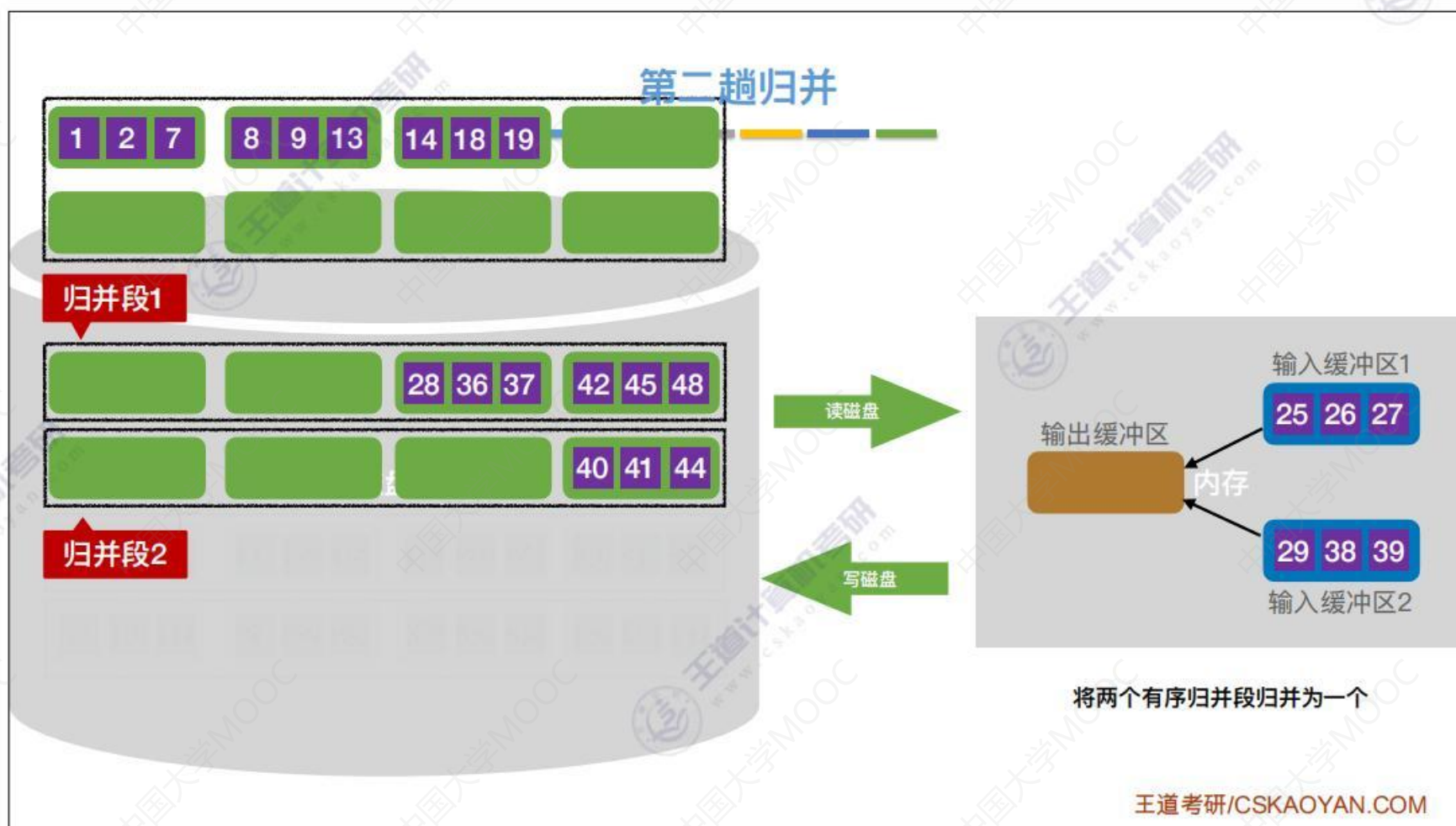
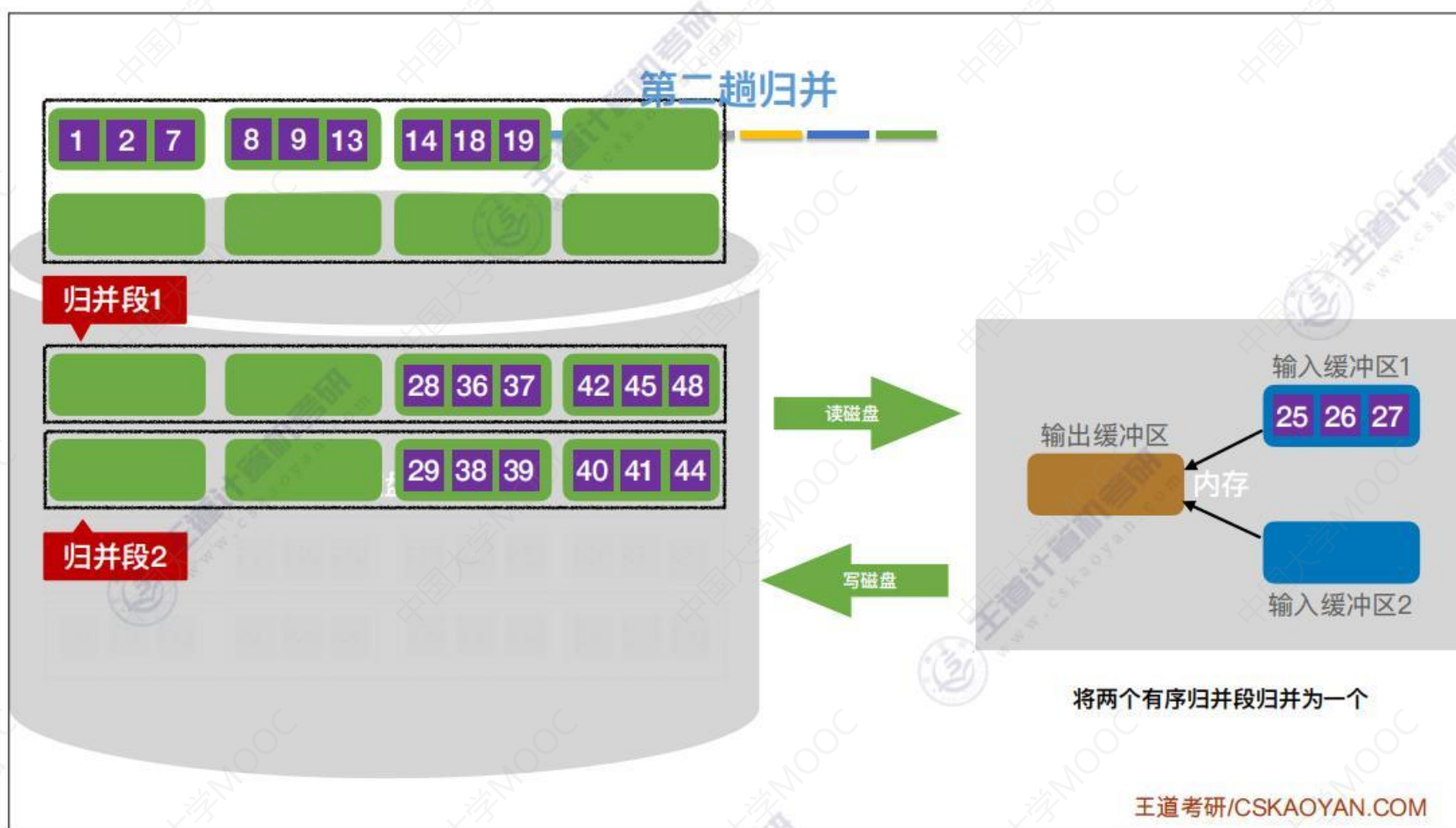


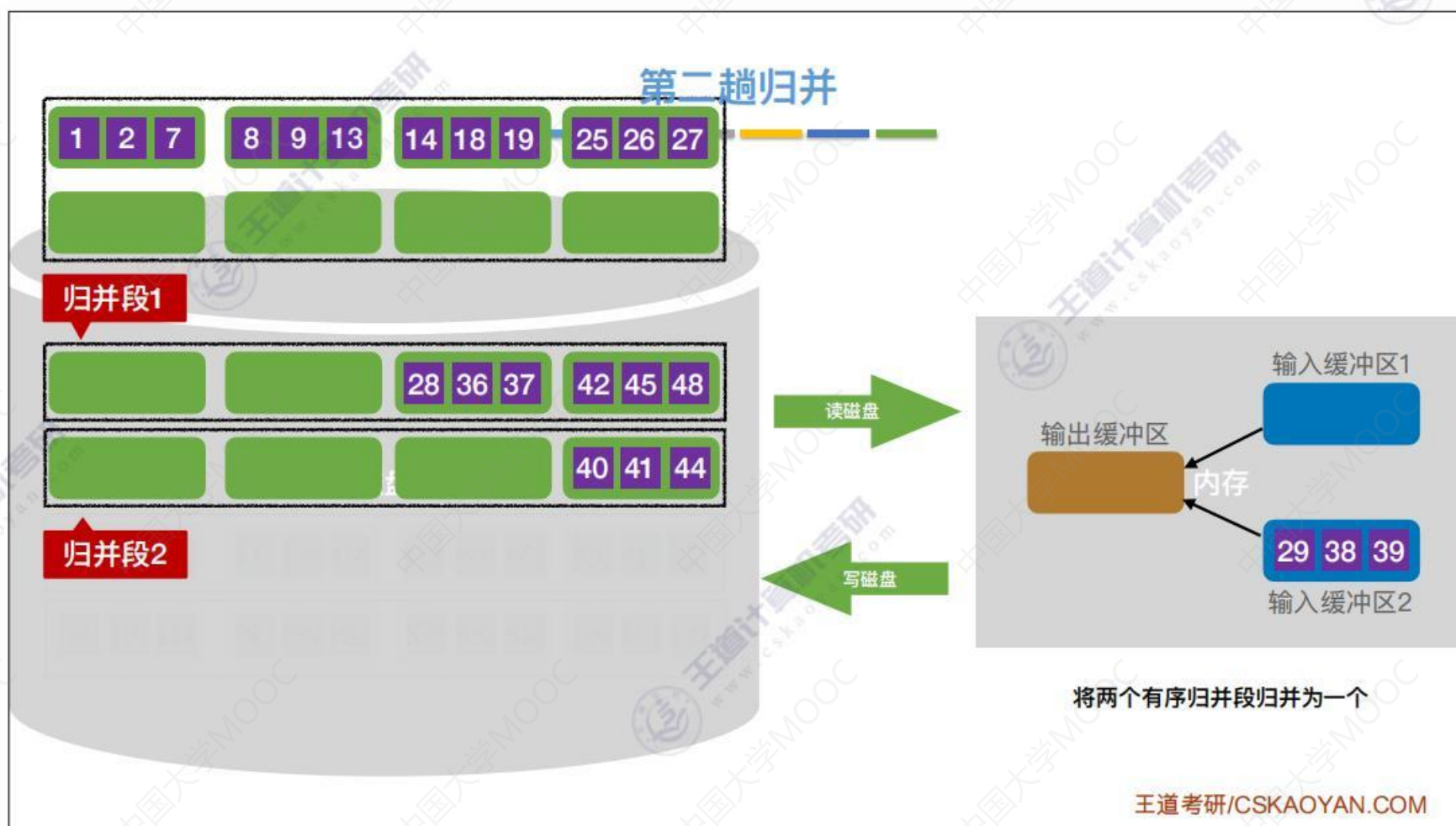
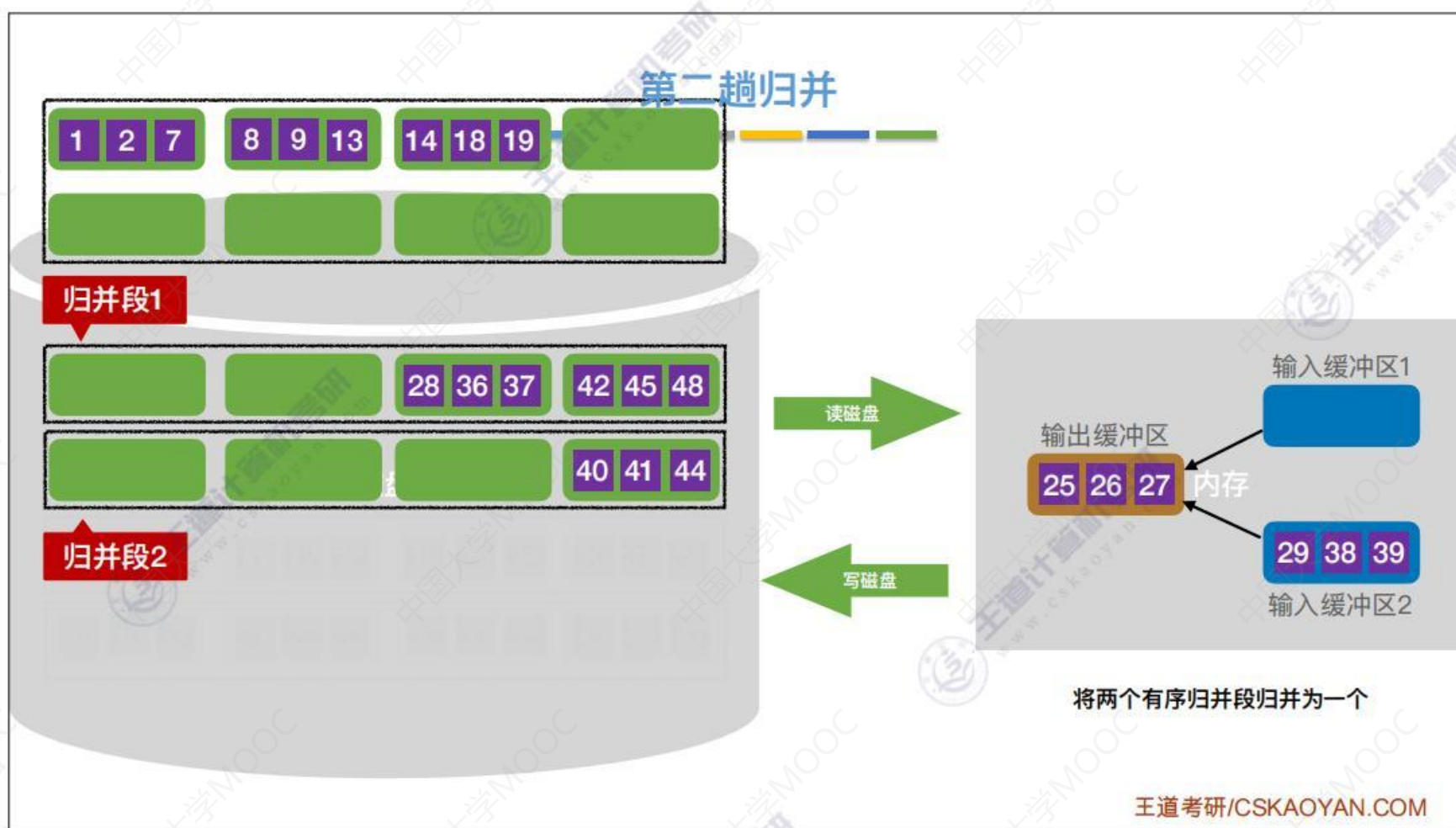


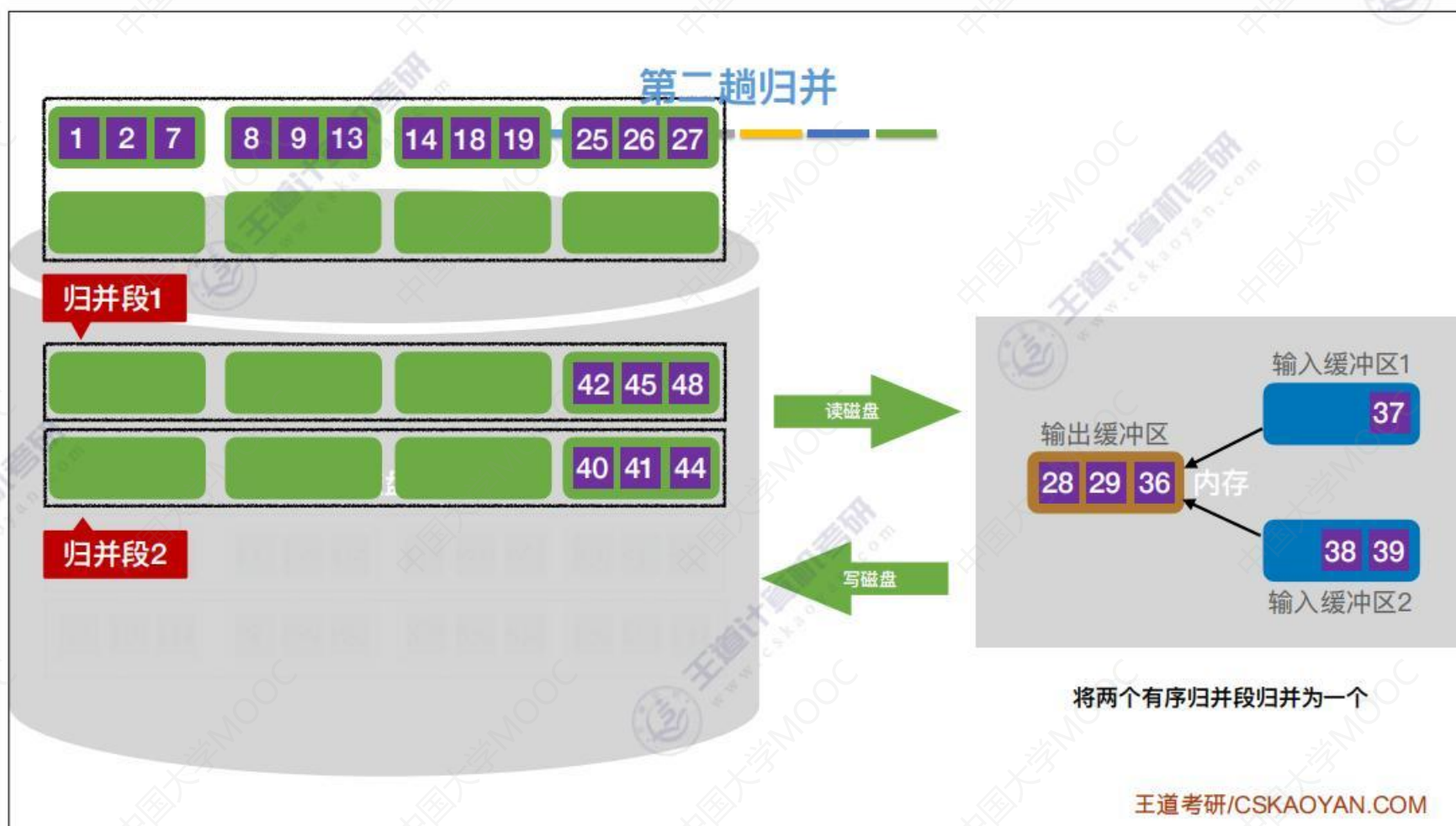
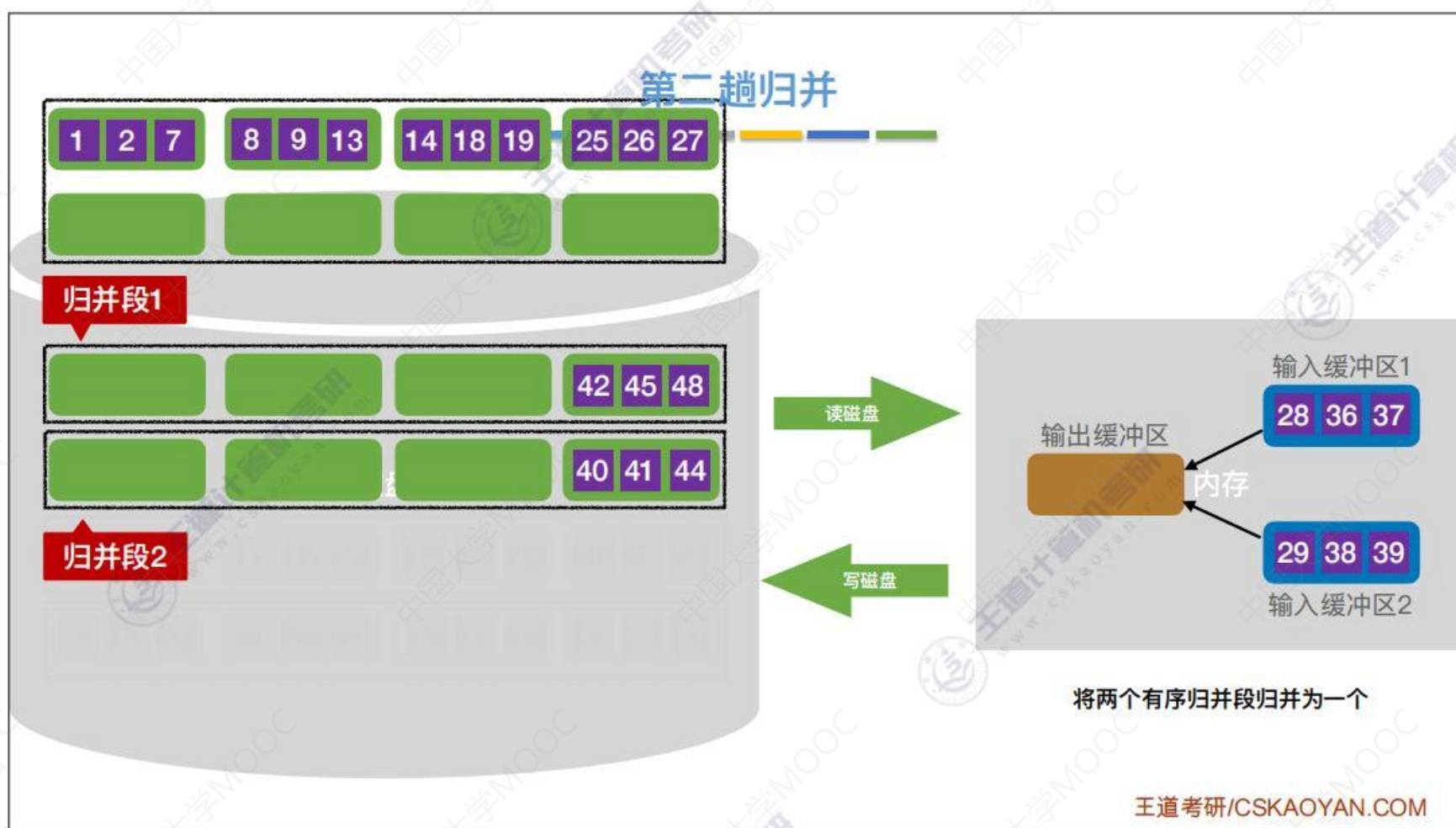


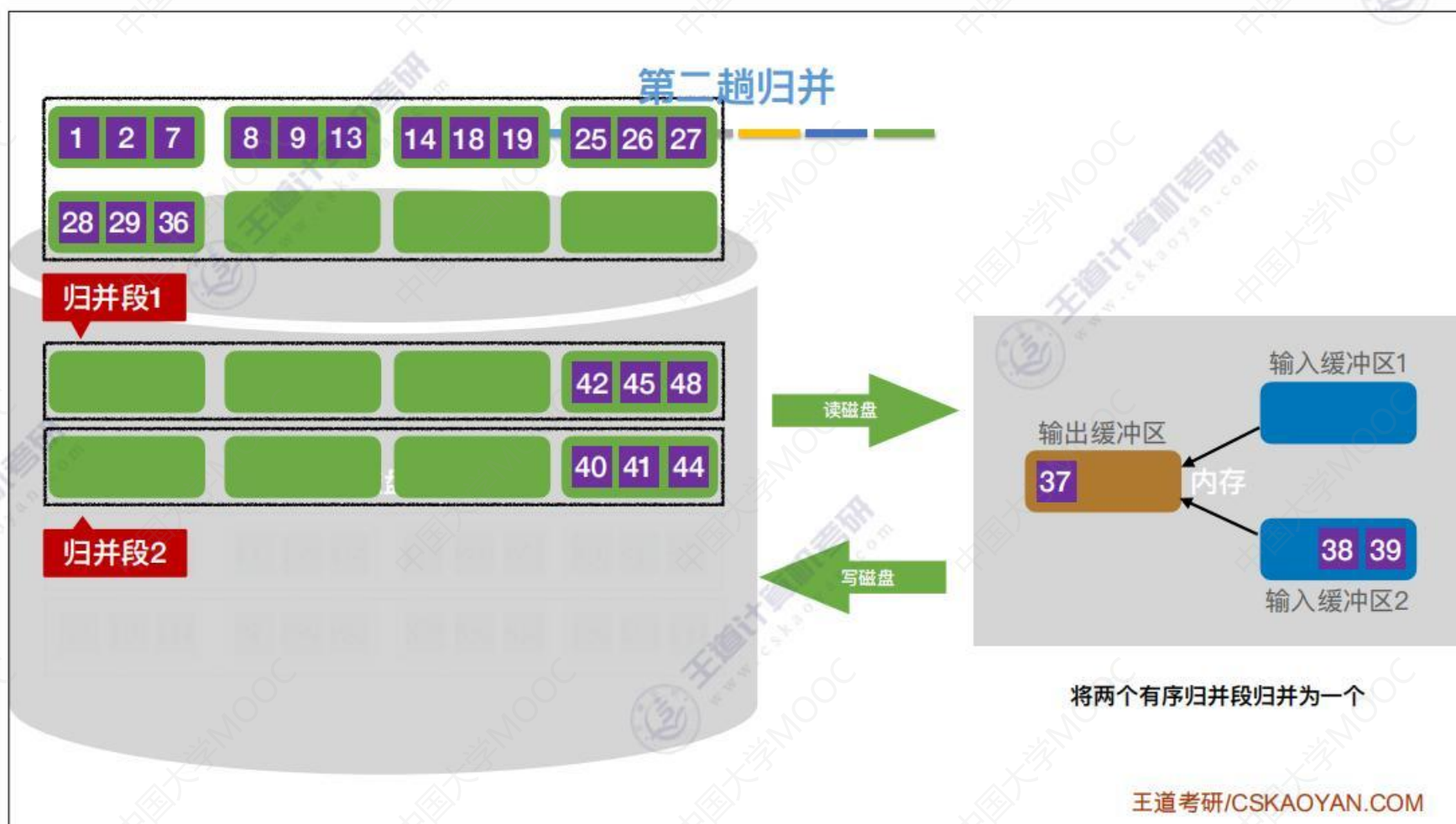
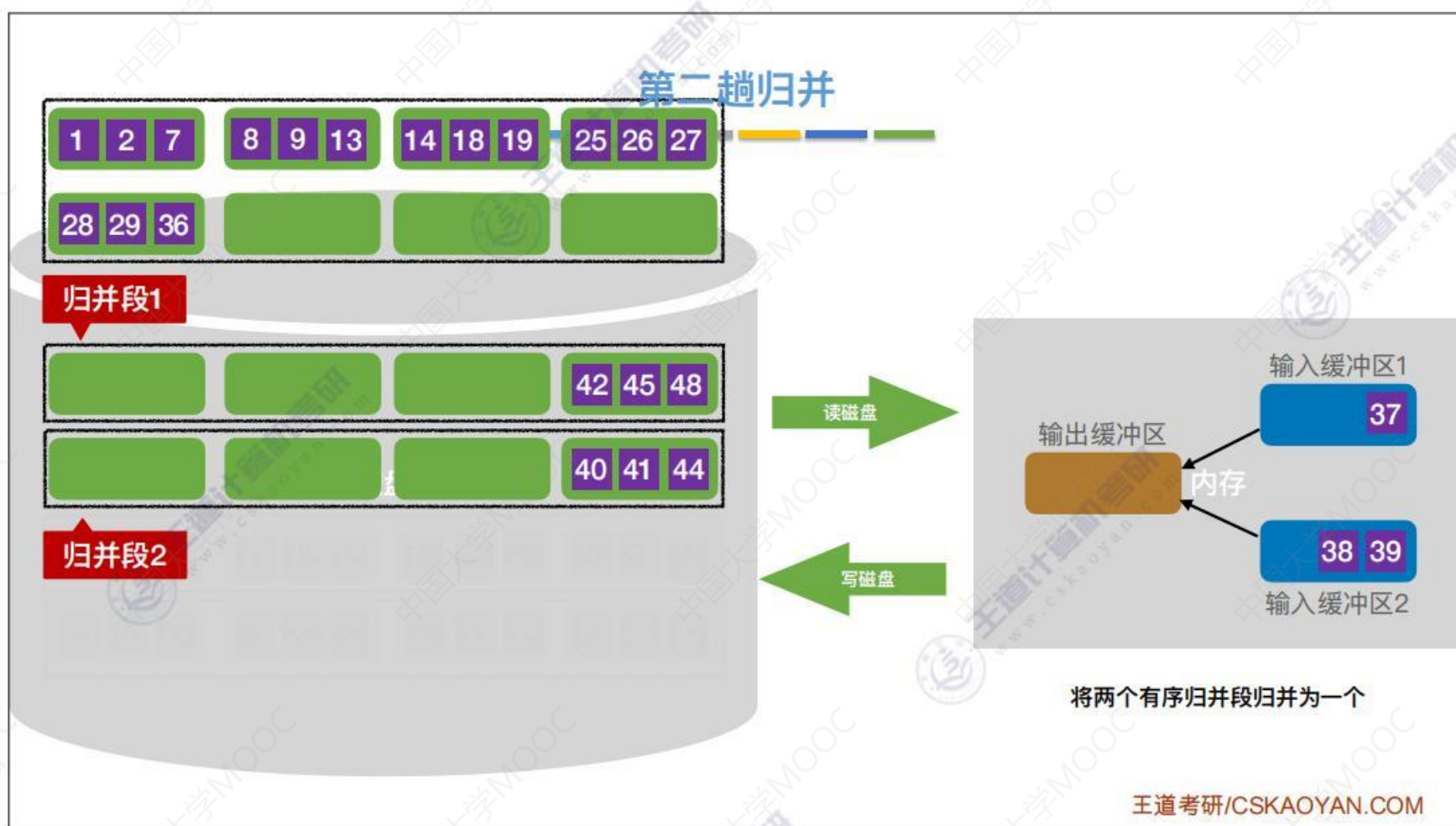


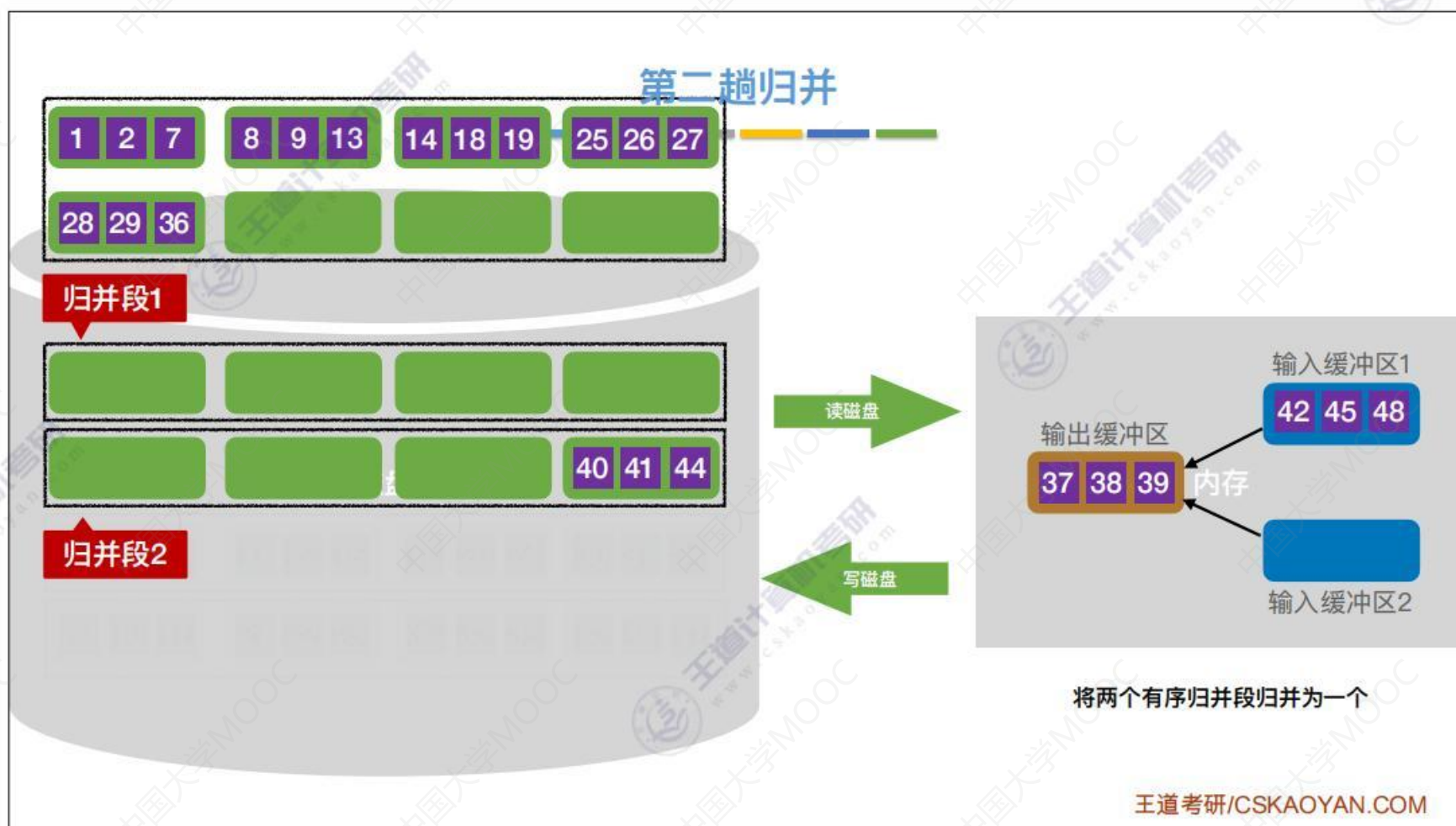
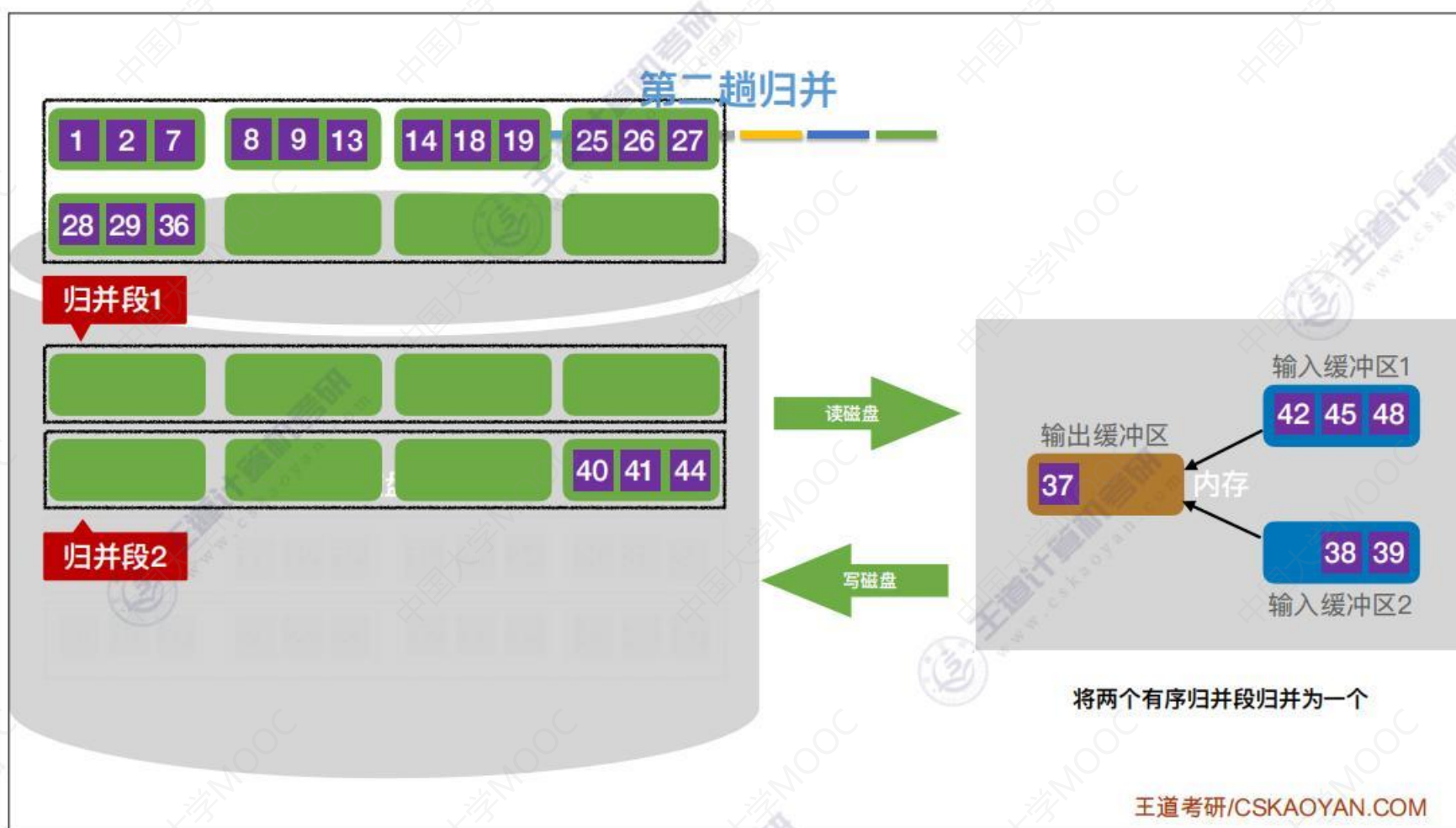


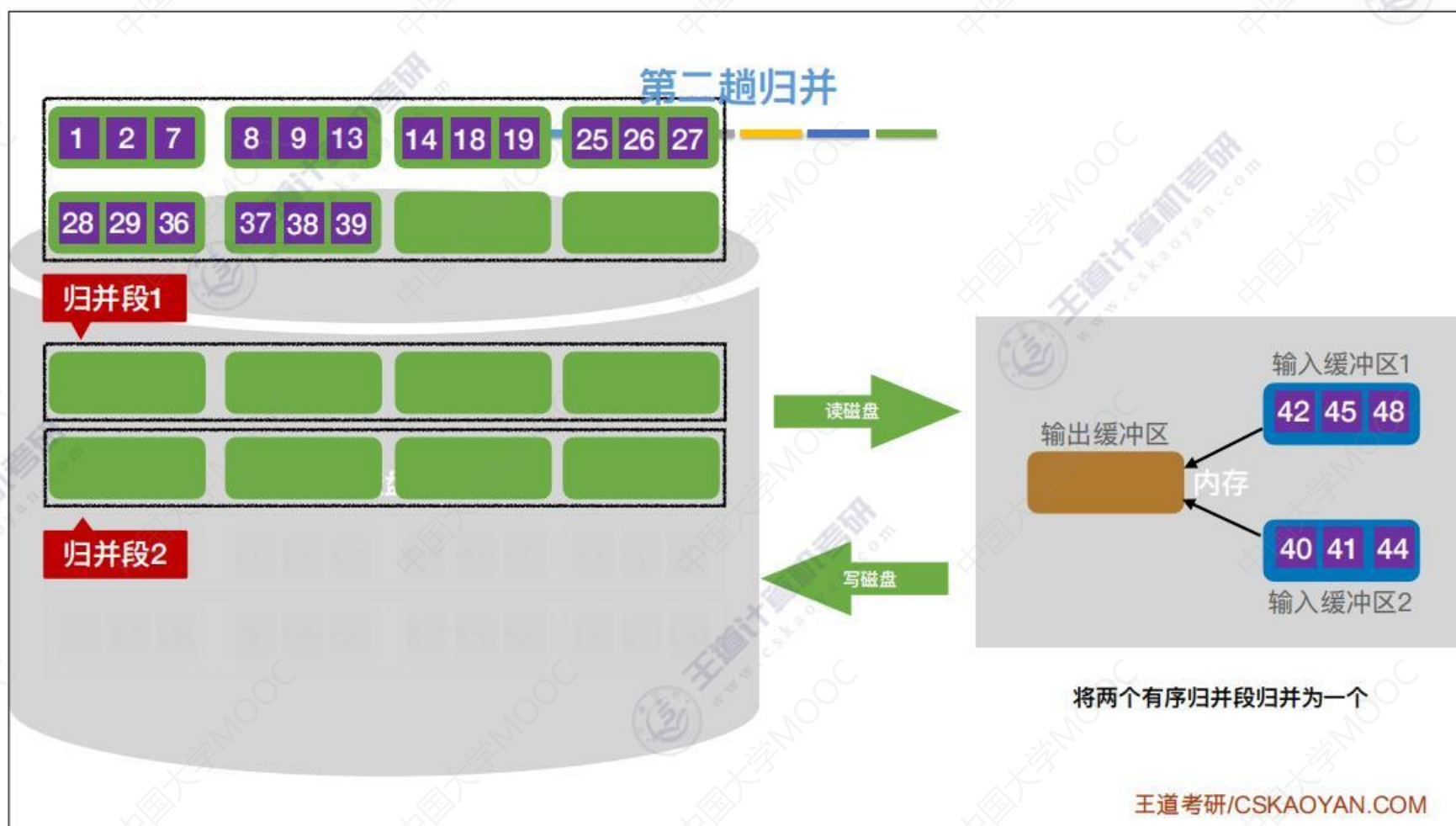
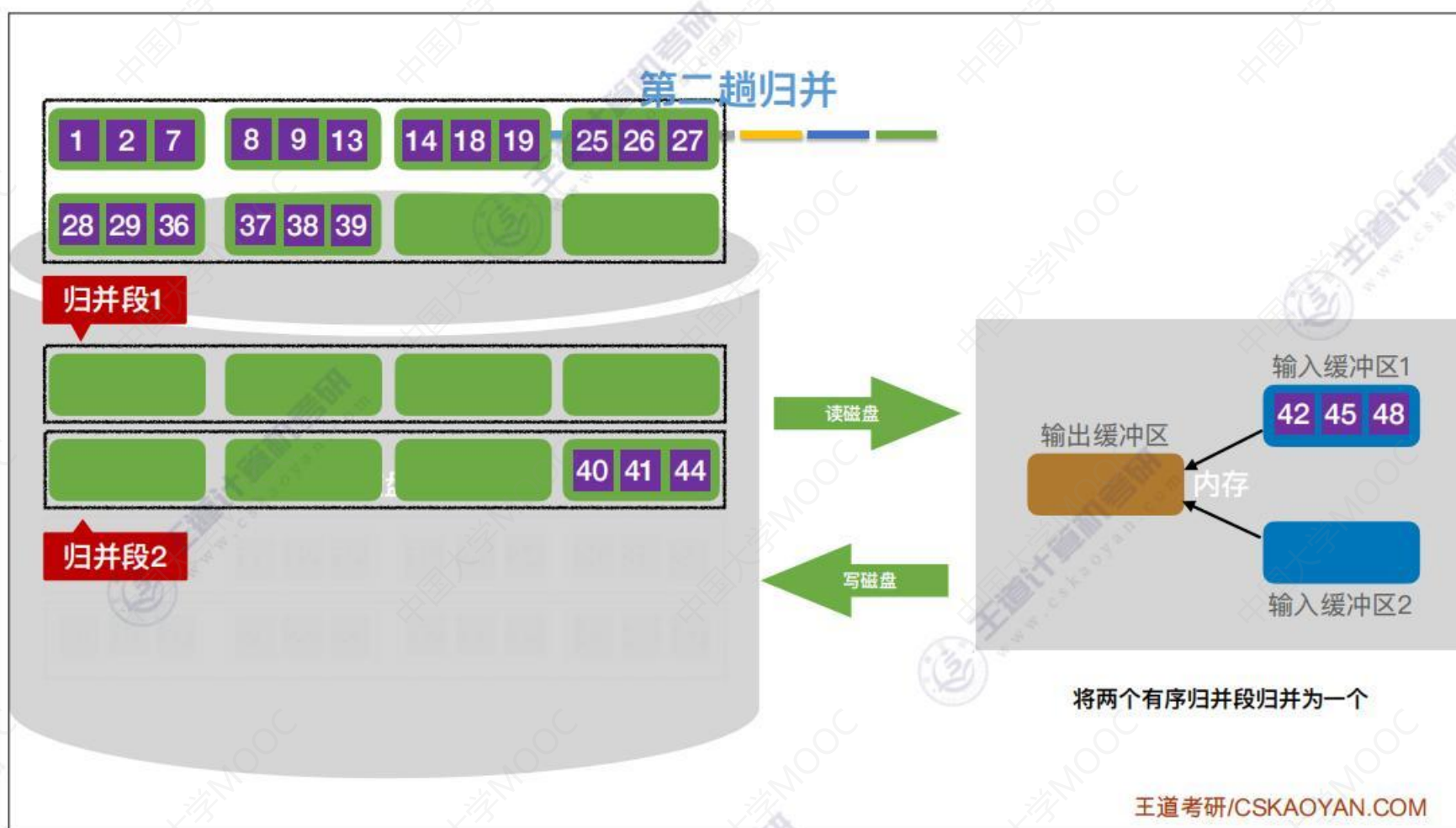


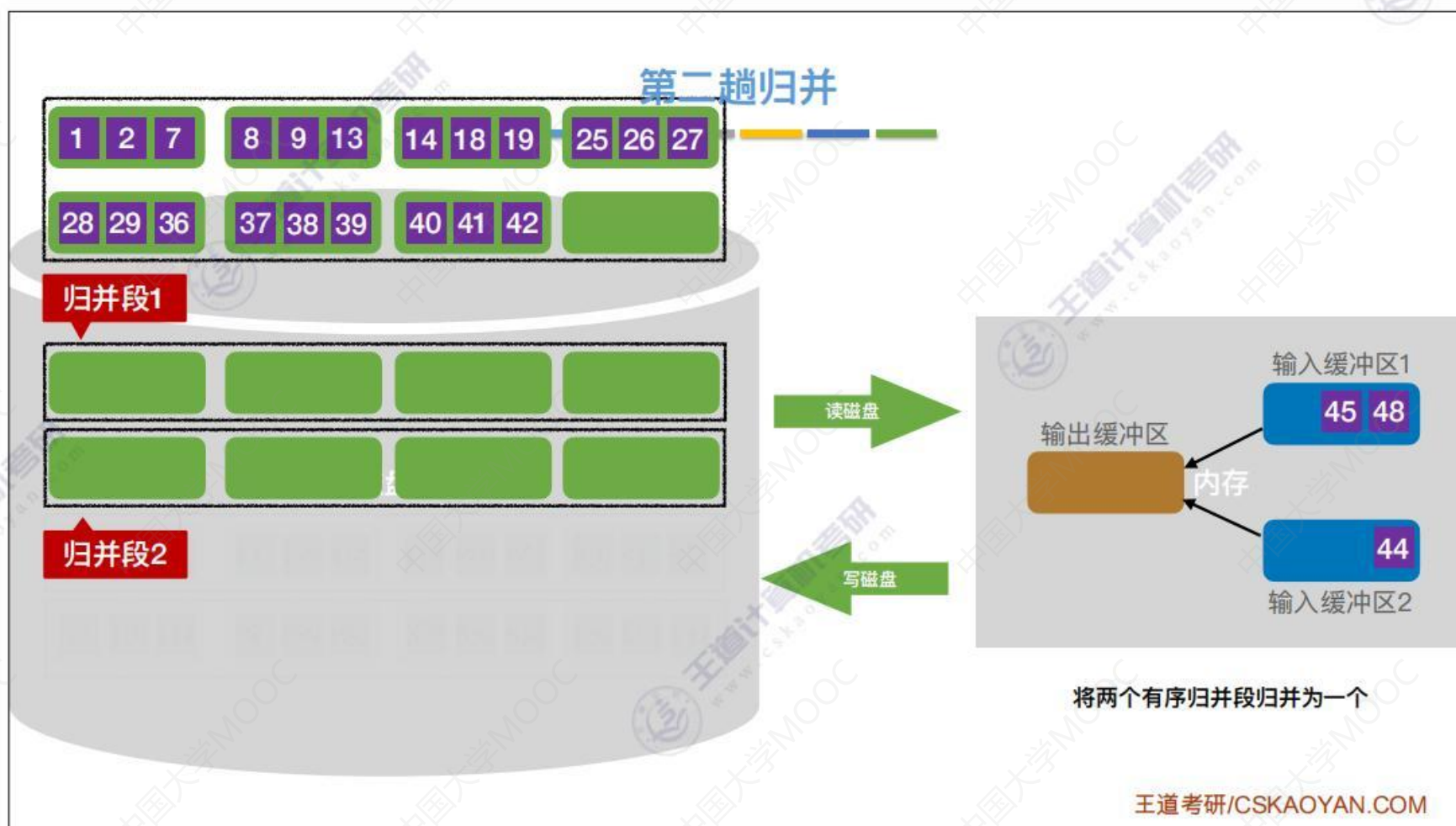
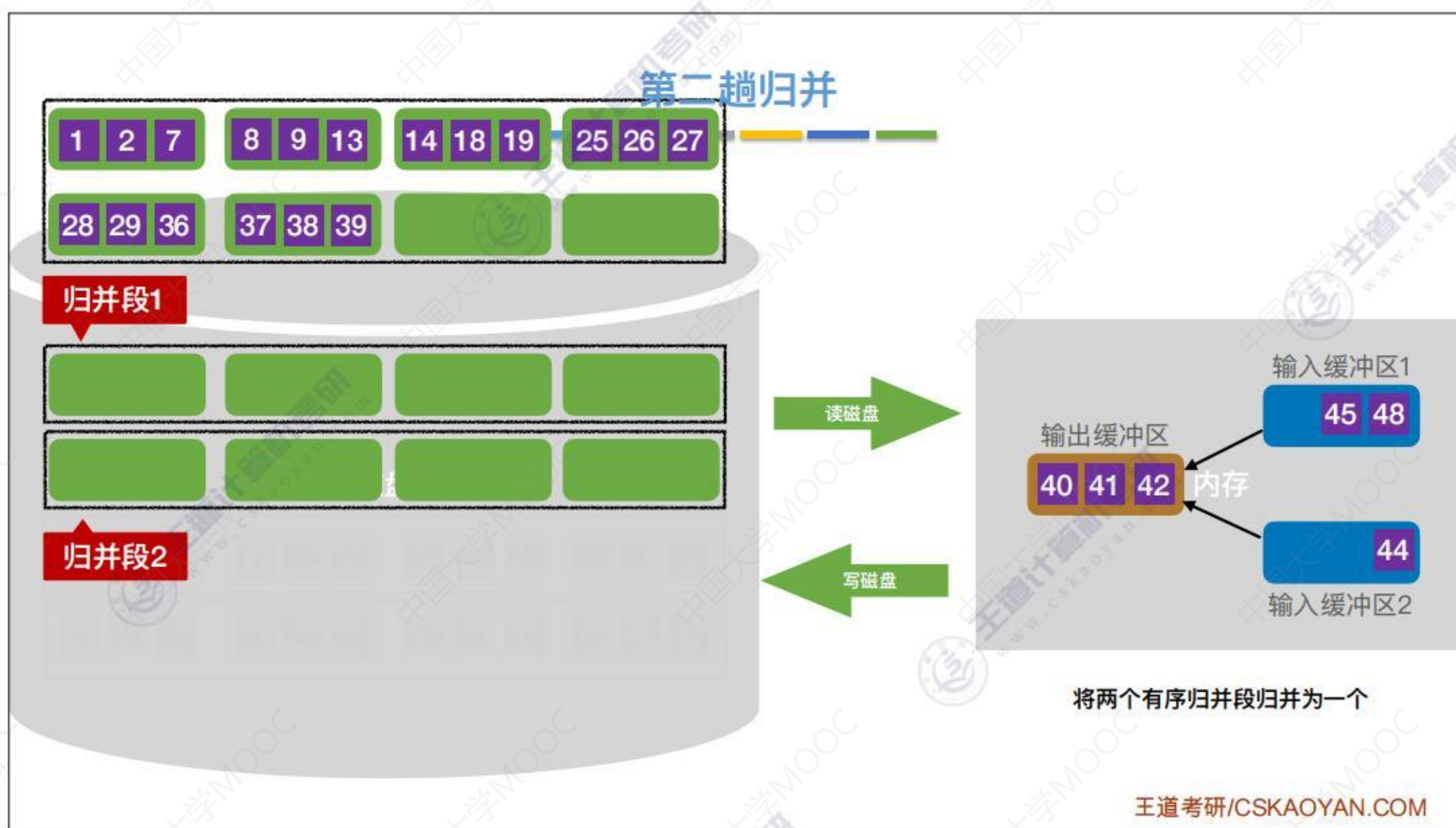


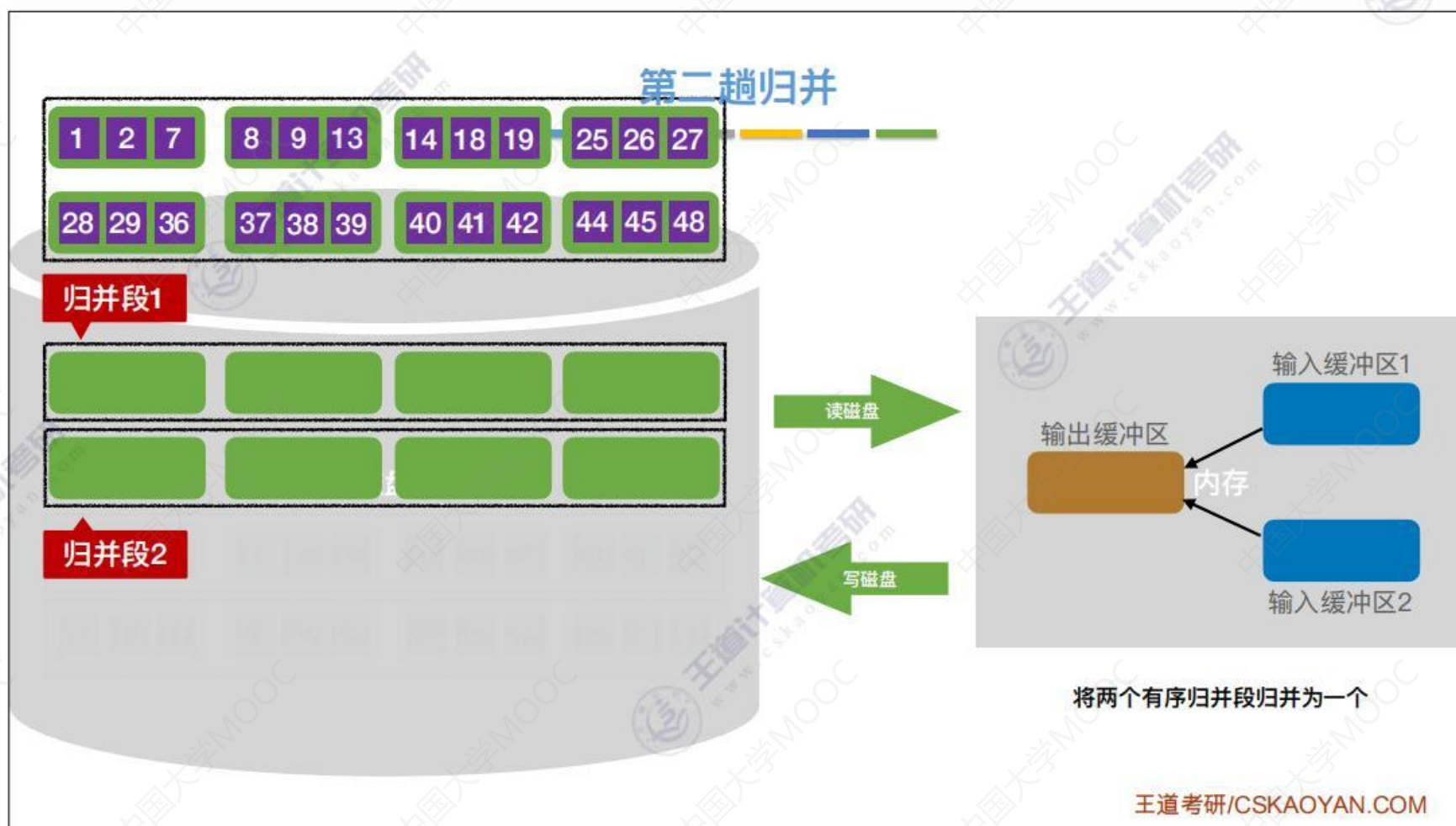
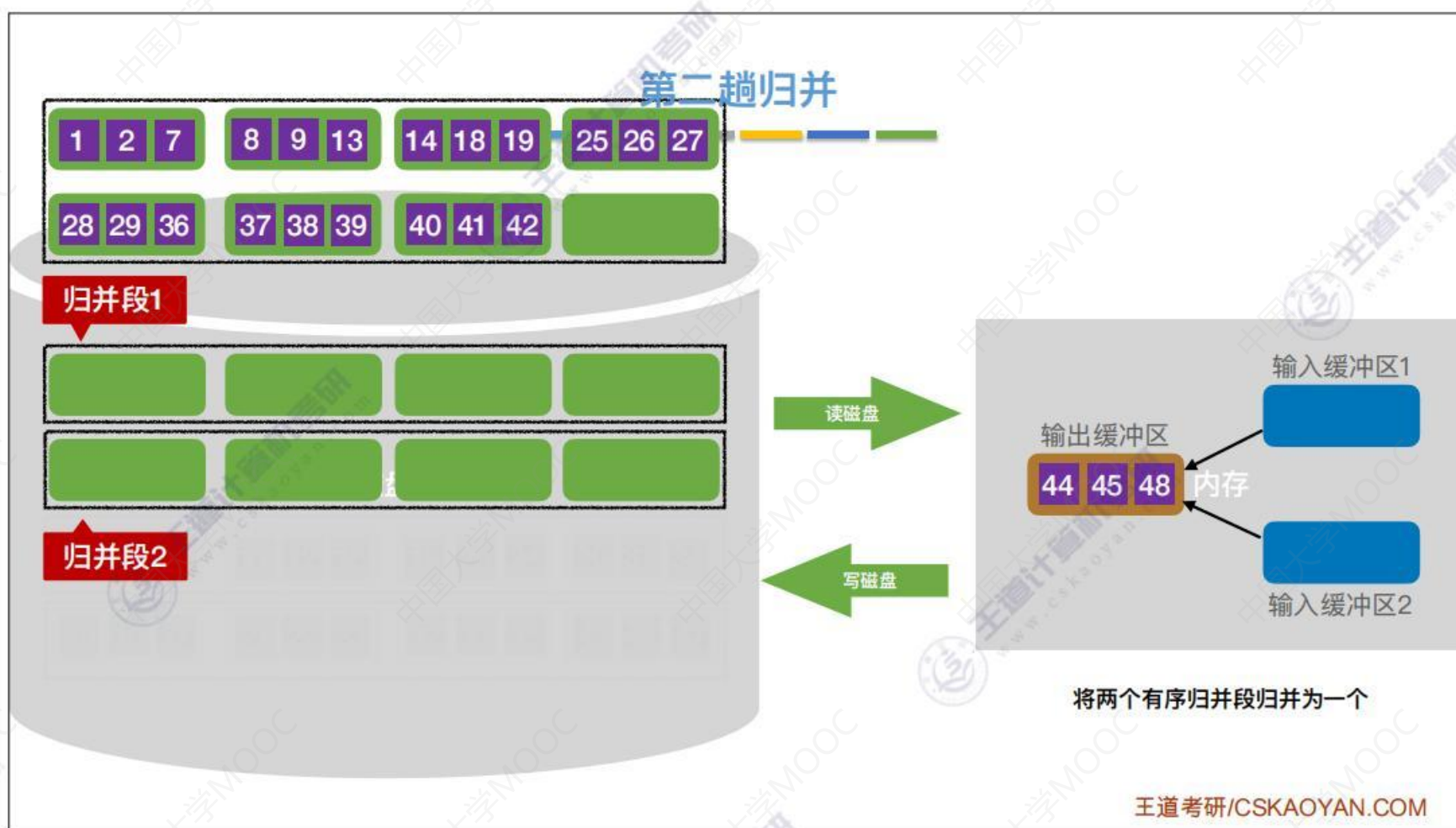


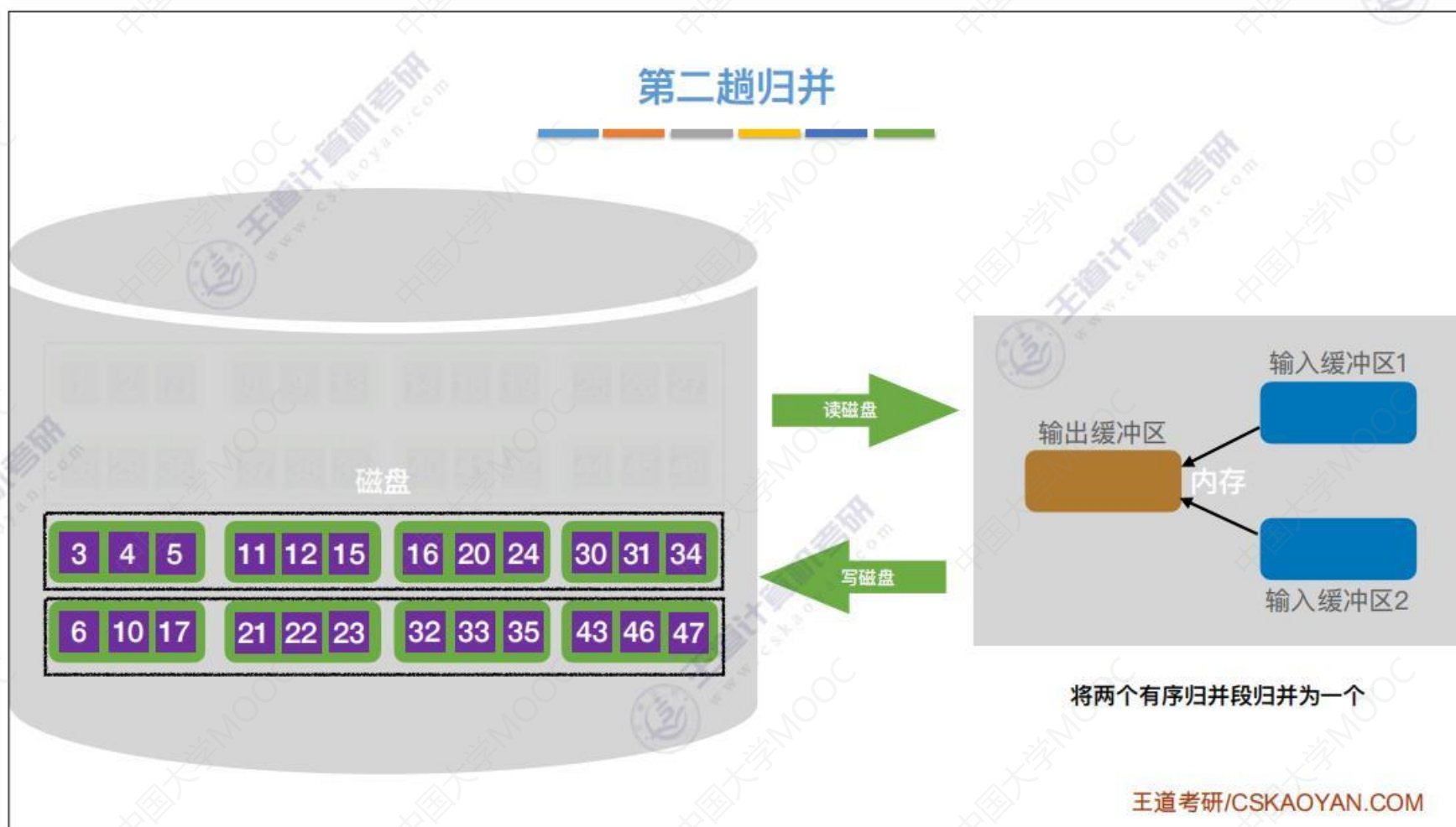
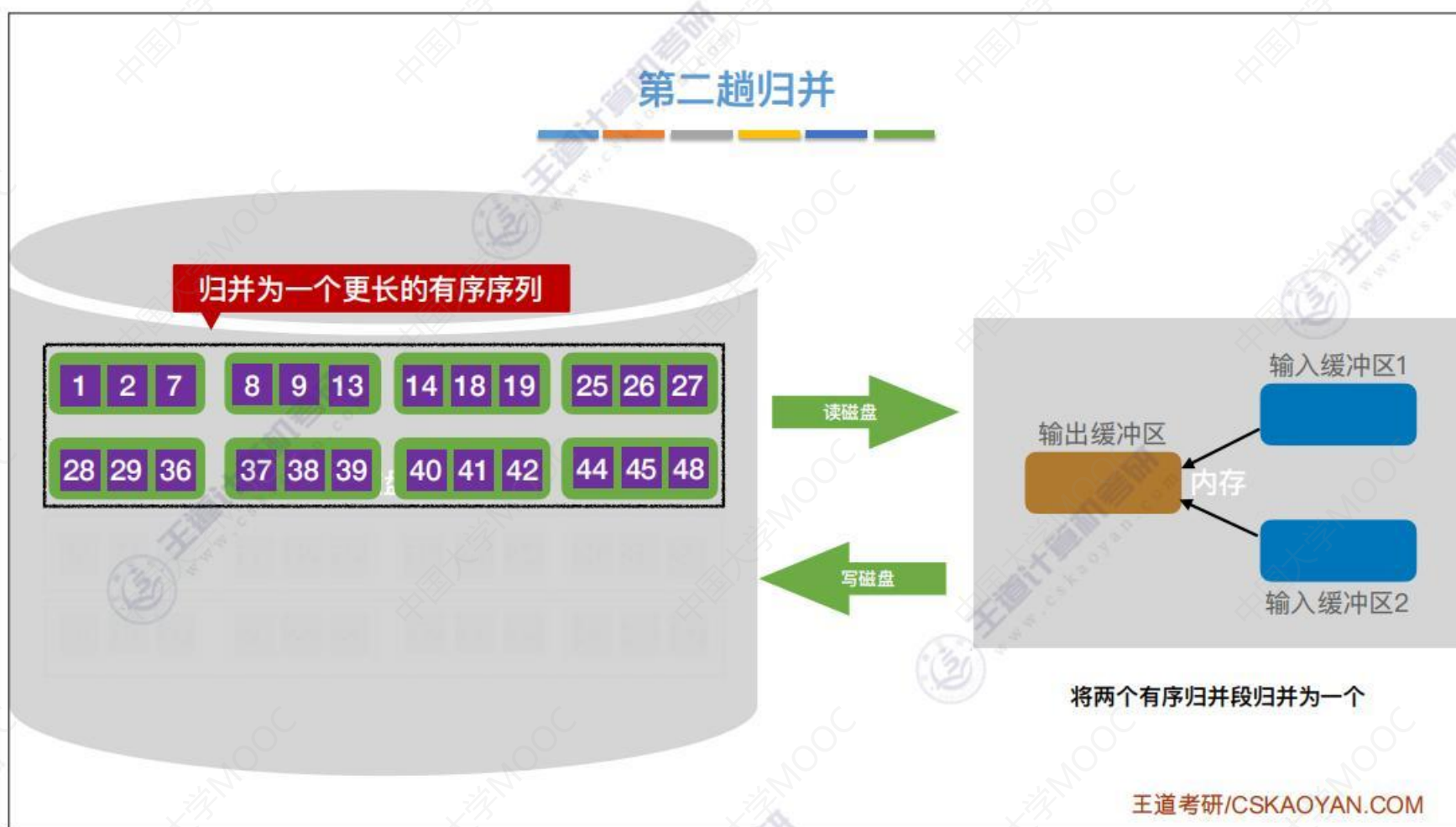


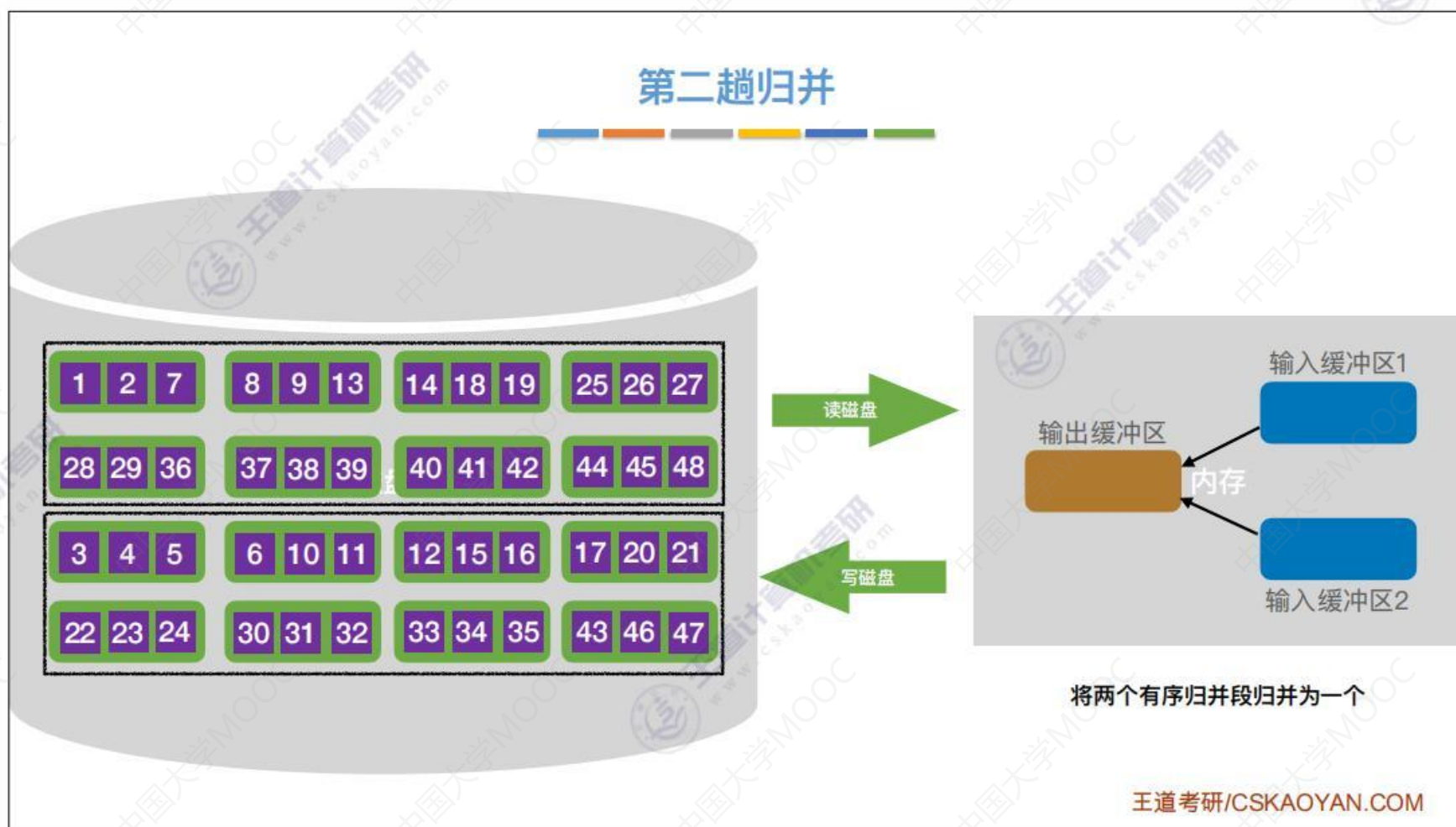
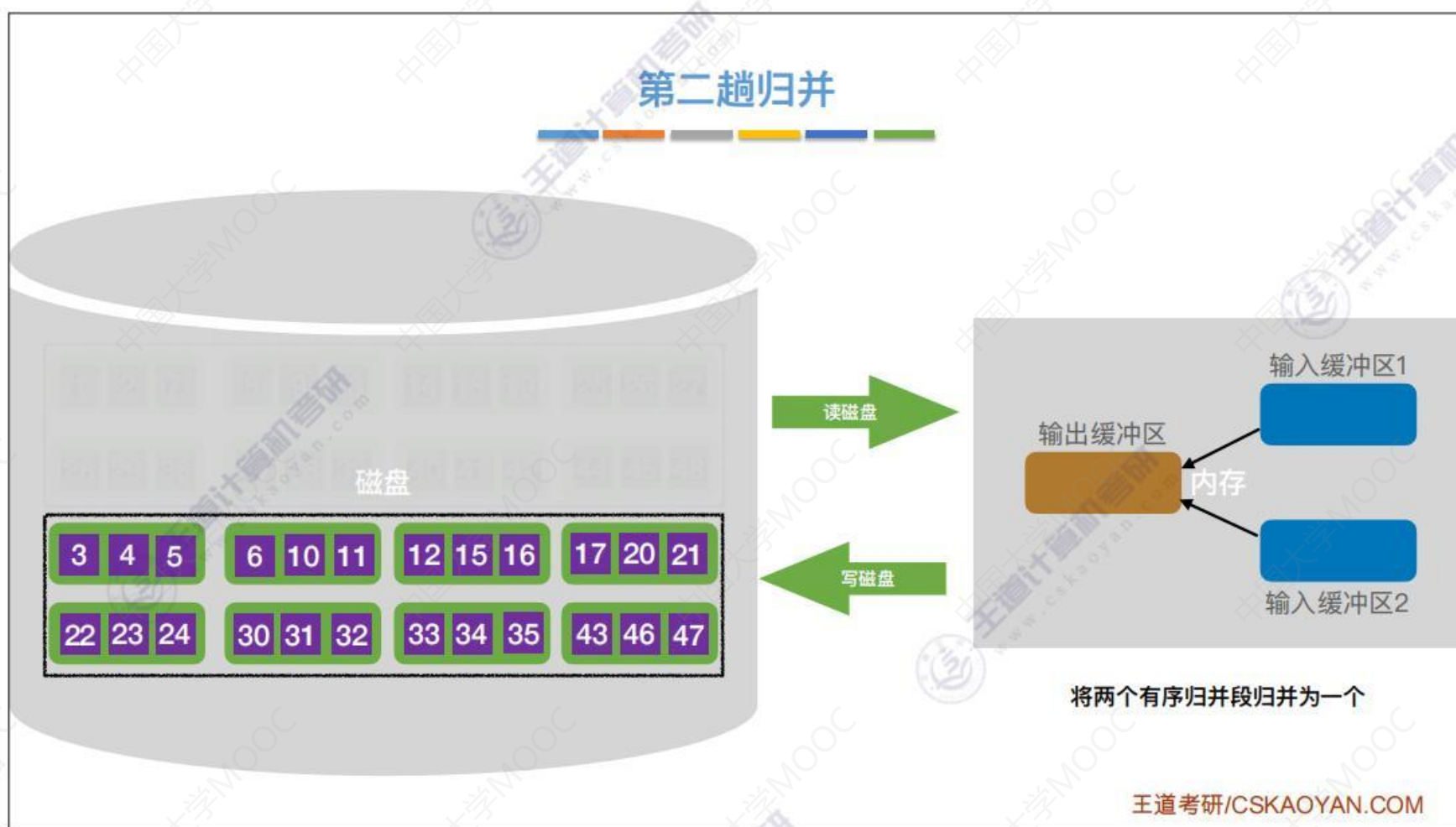


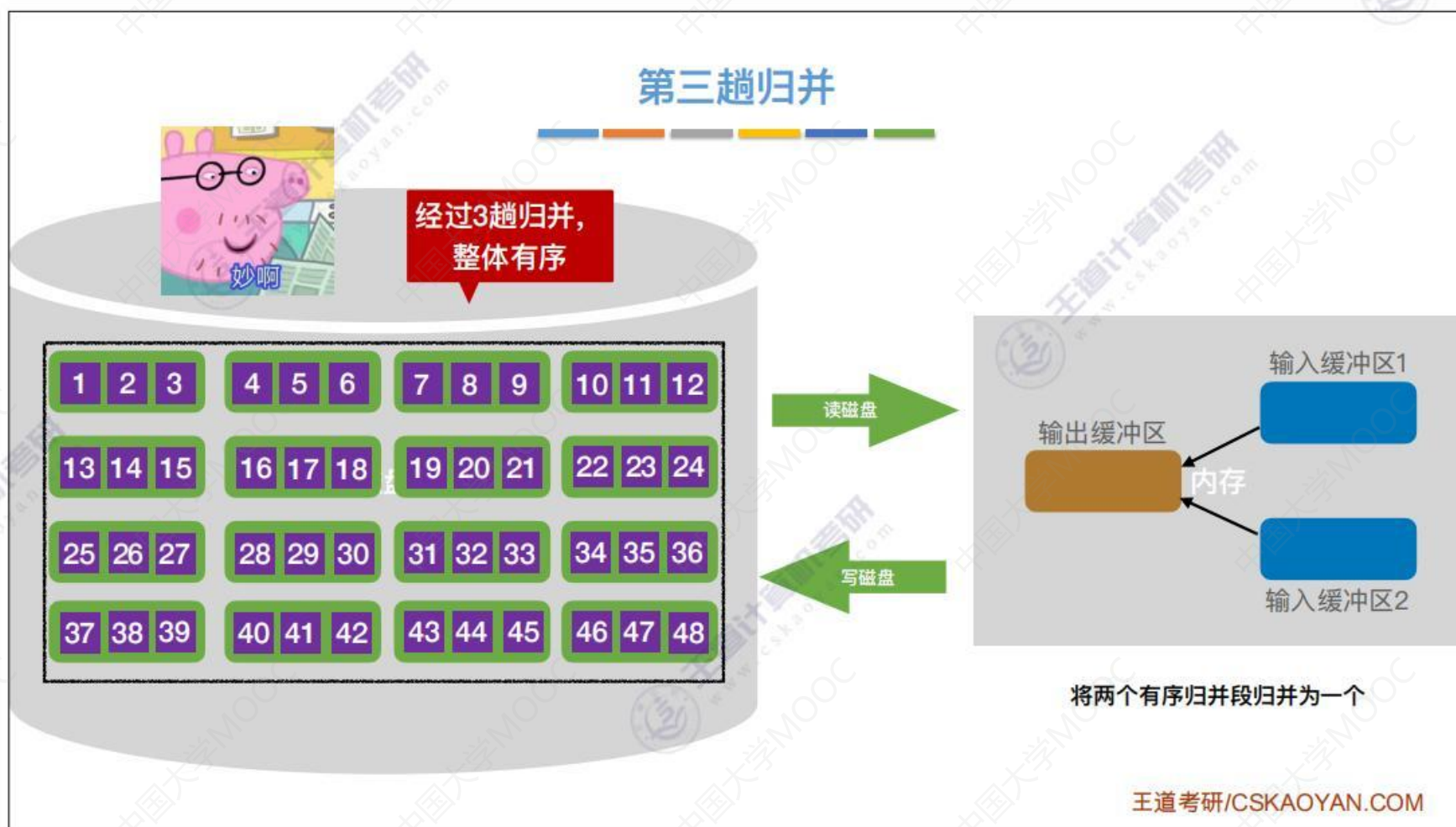




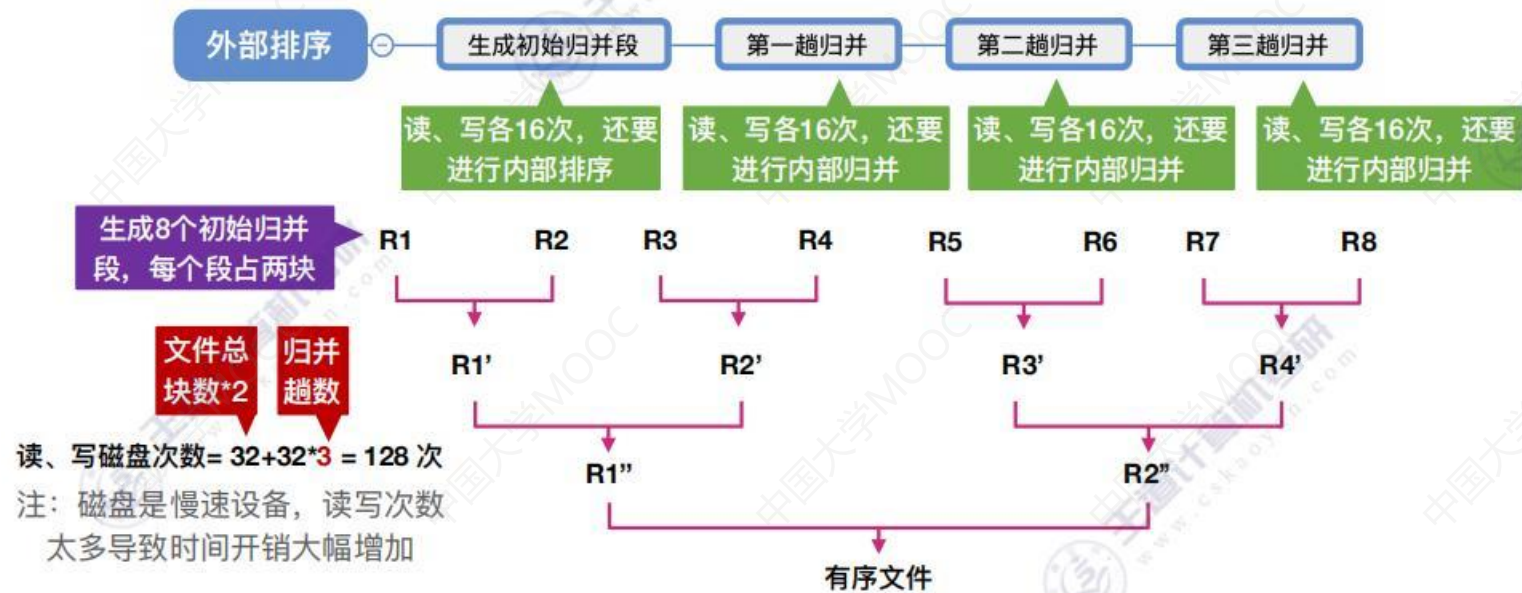








## 时间开销分析



外部排序时间开销 = 读写外存的时间 + 内部排序所需时间 + 内部归并所需时间

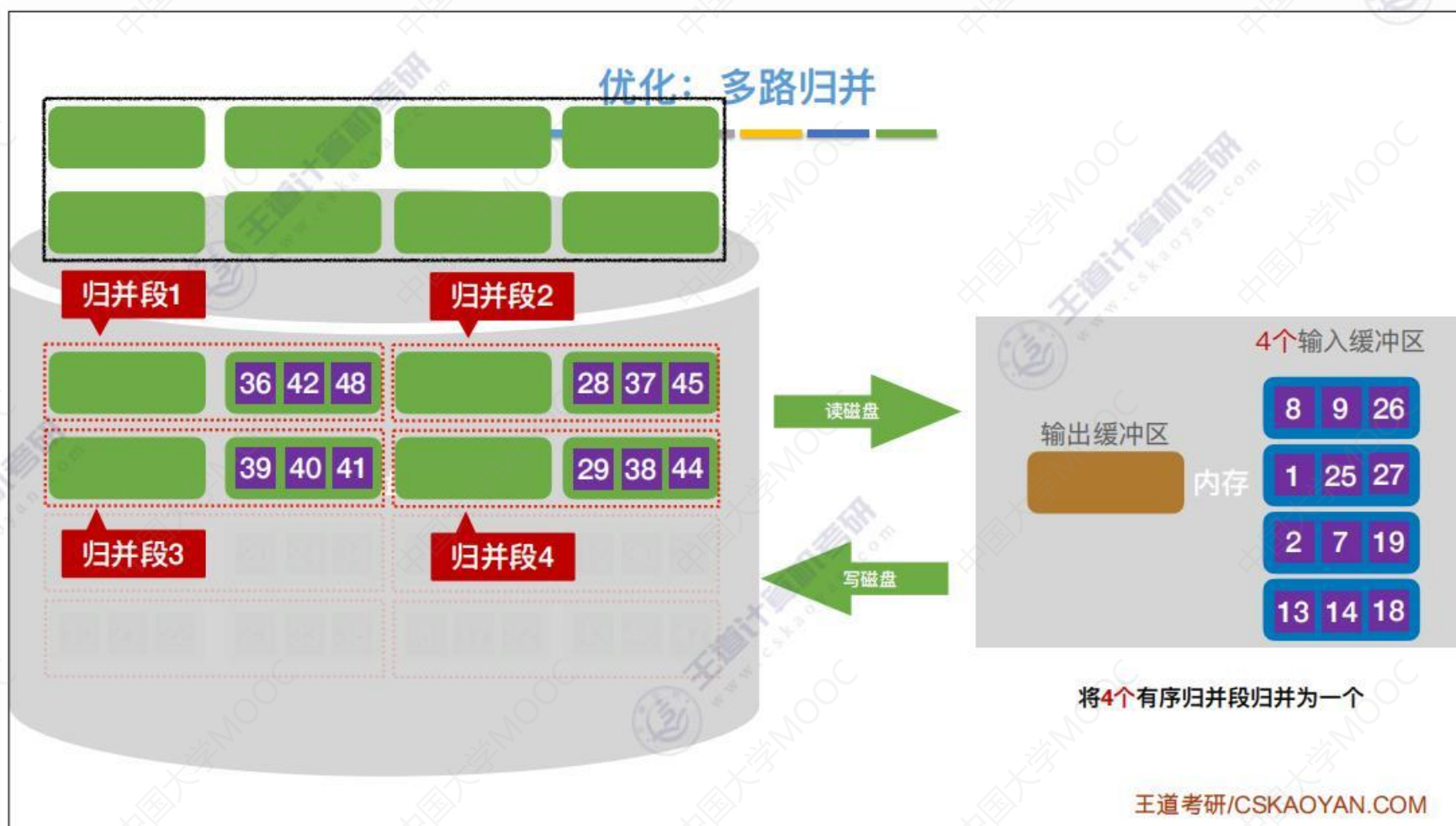
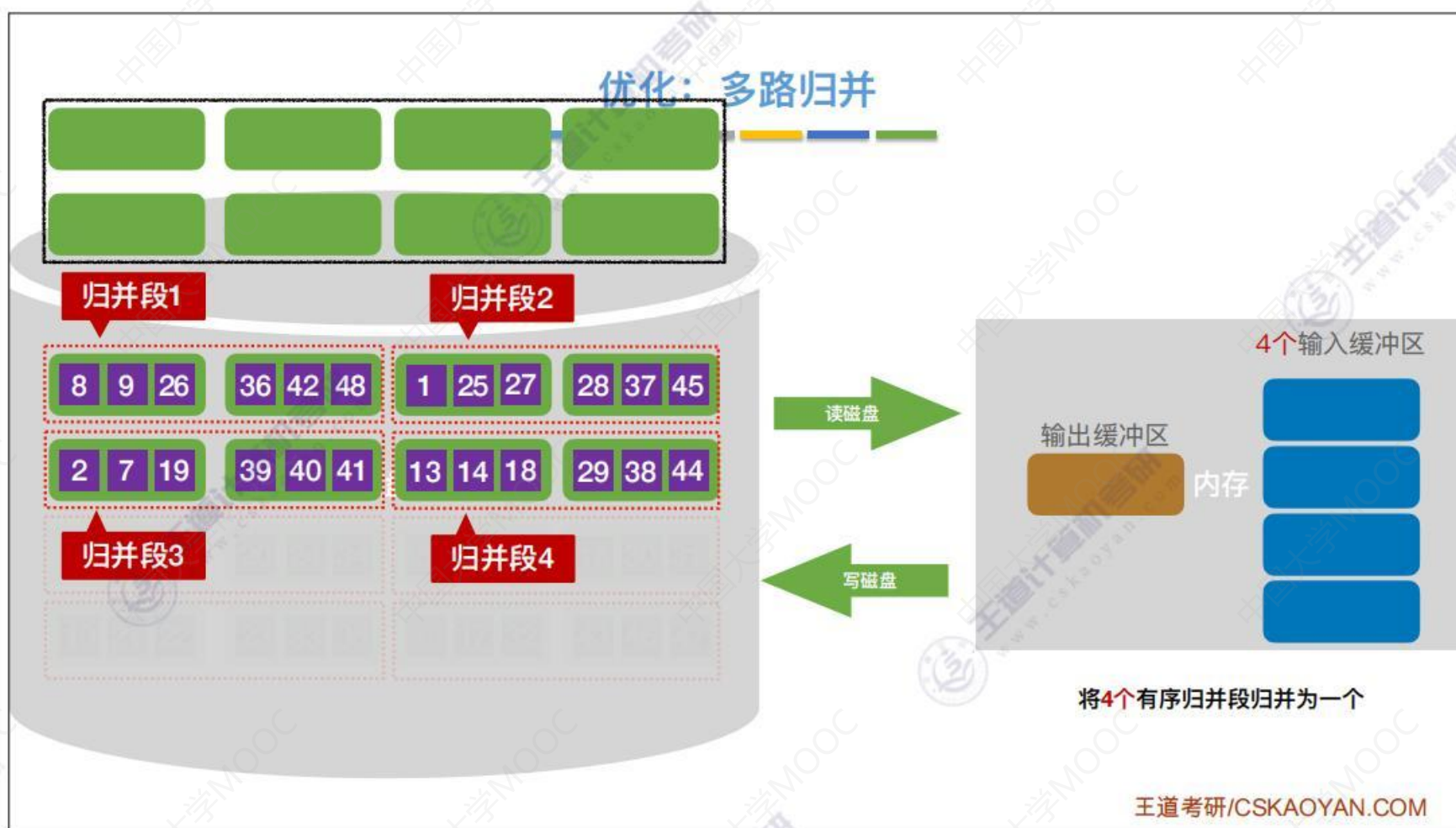
王道考研/CSKAOYAN.COM

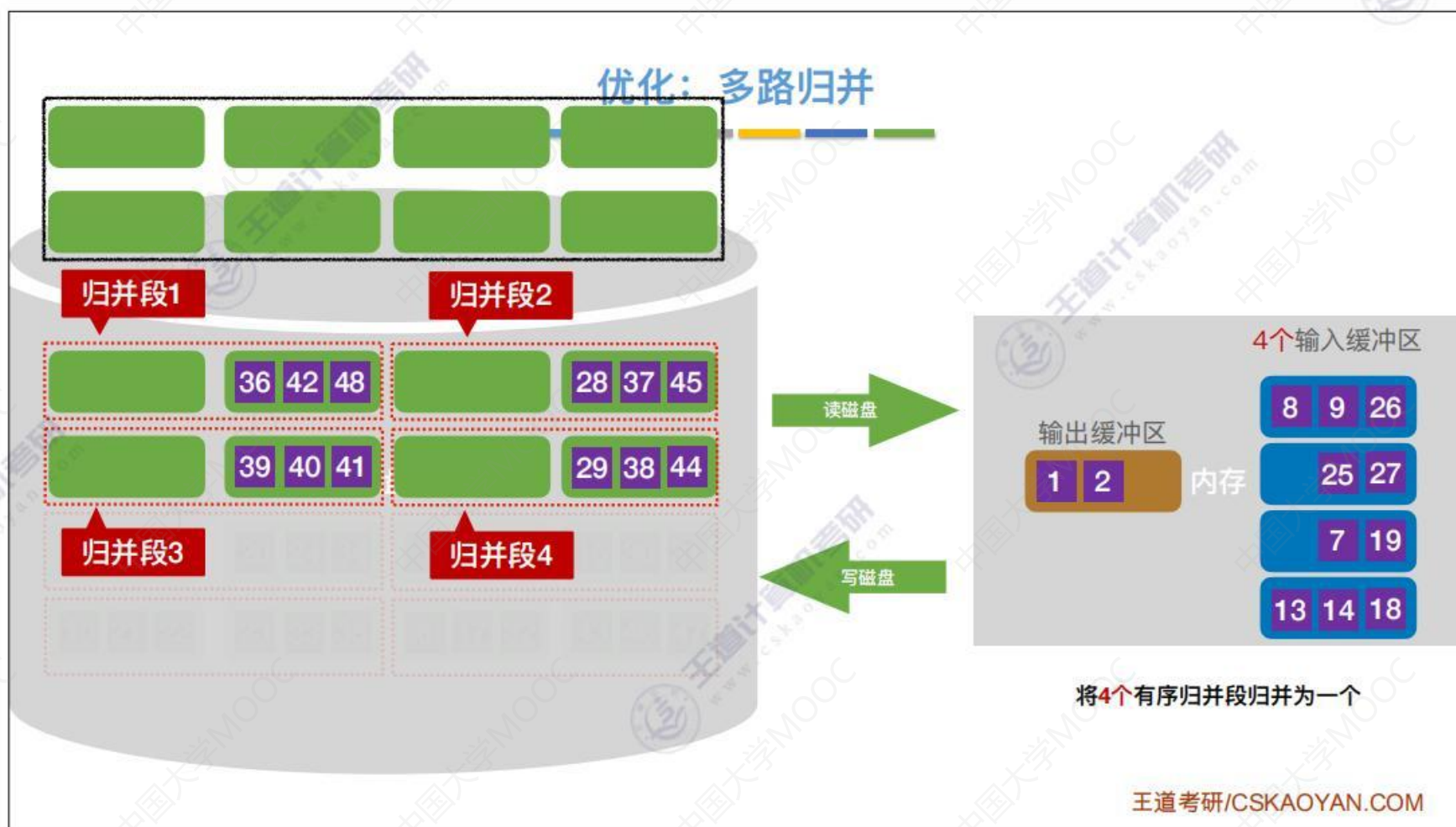
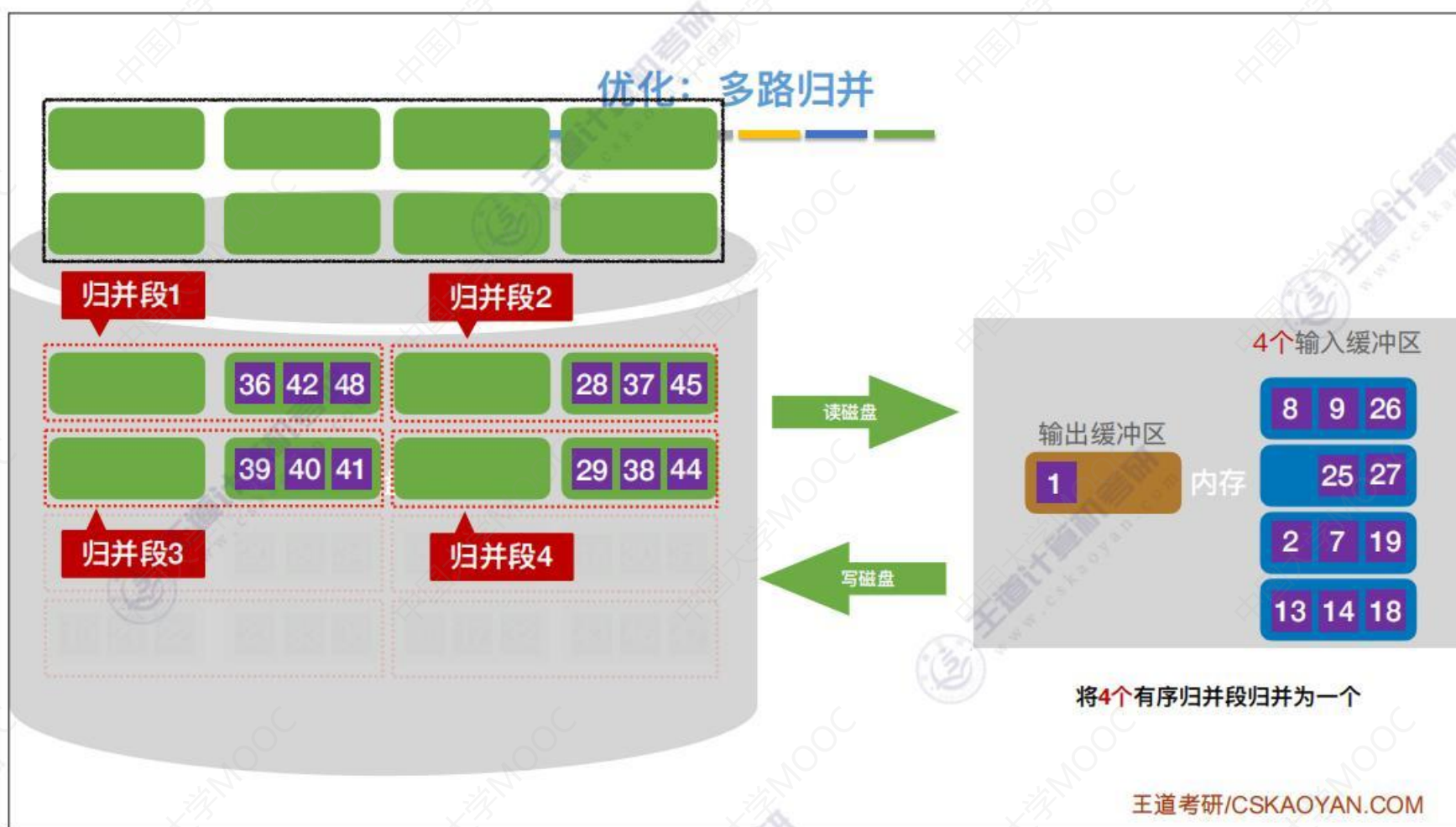
## 如何优化？

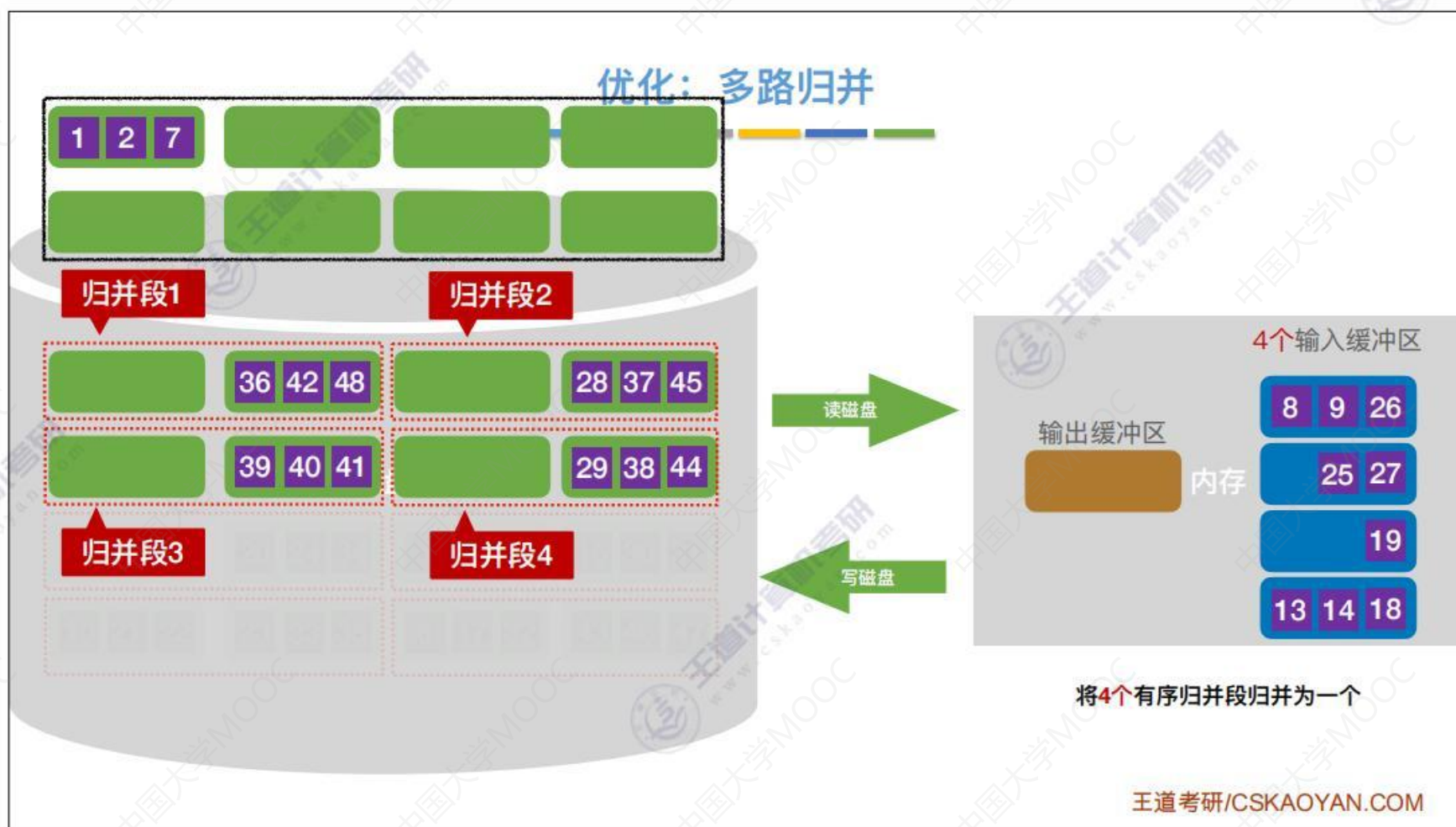
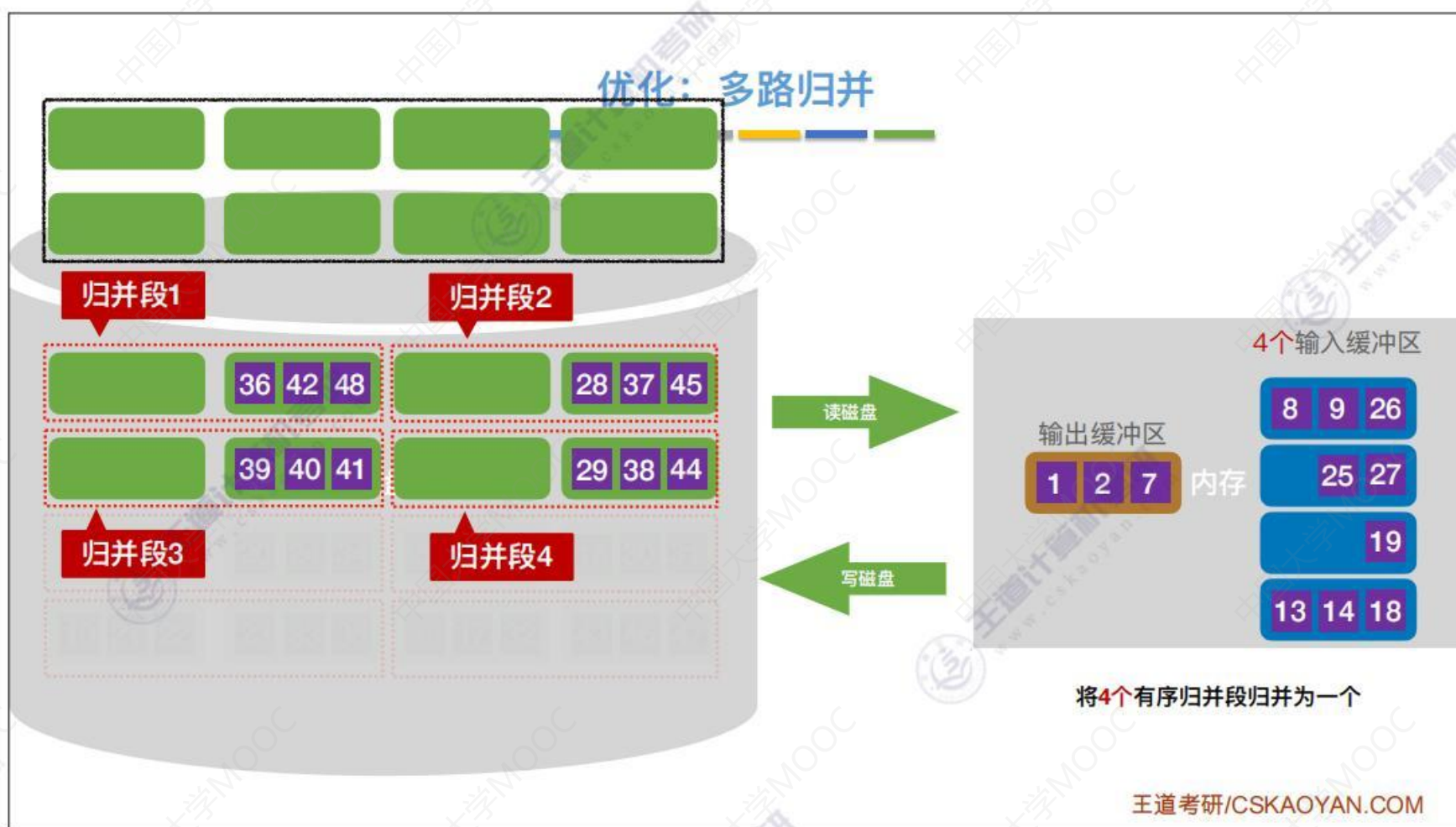
外部排序时间开销 = 读写外存的时间 + 内部排序所需时间 + 内部归并所需时间

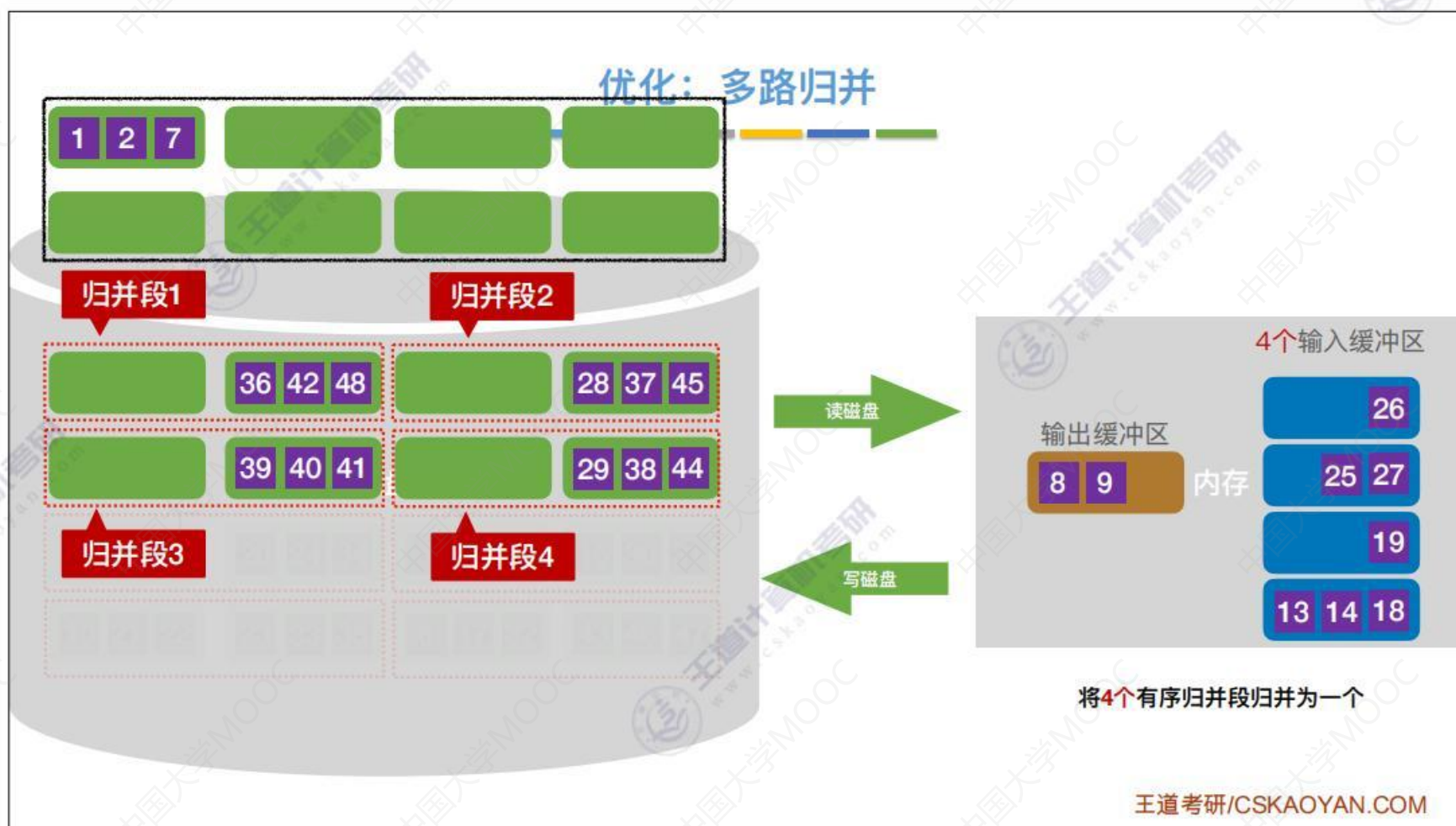
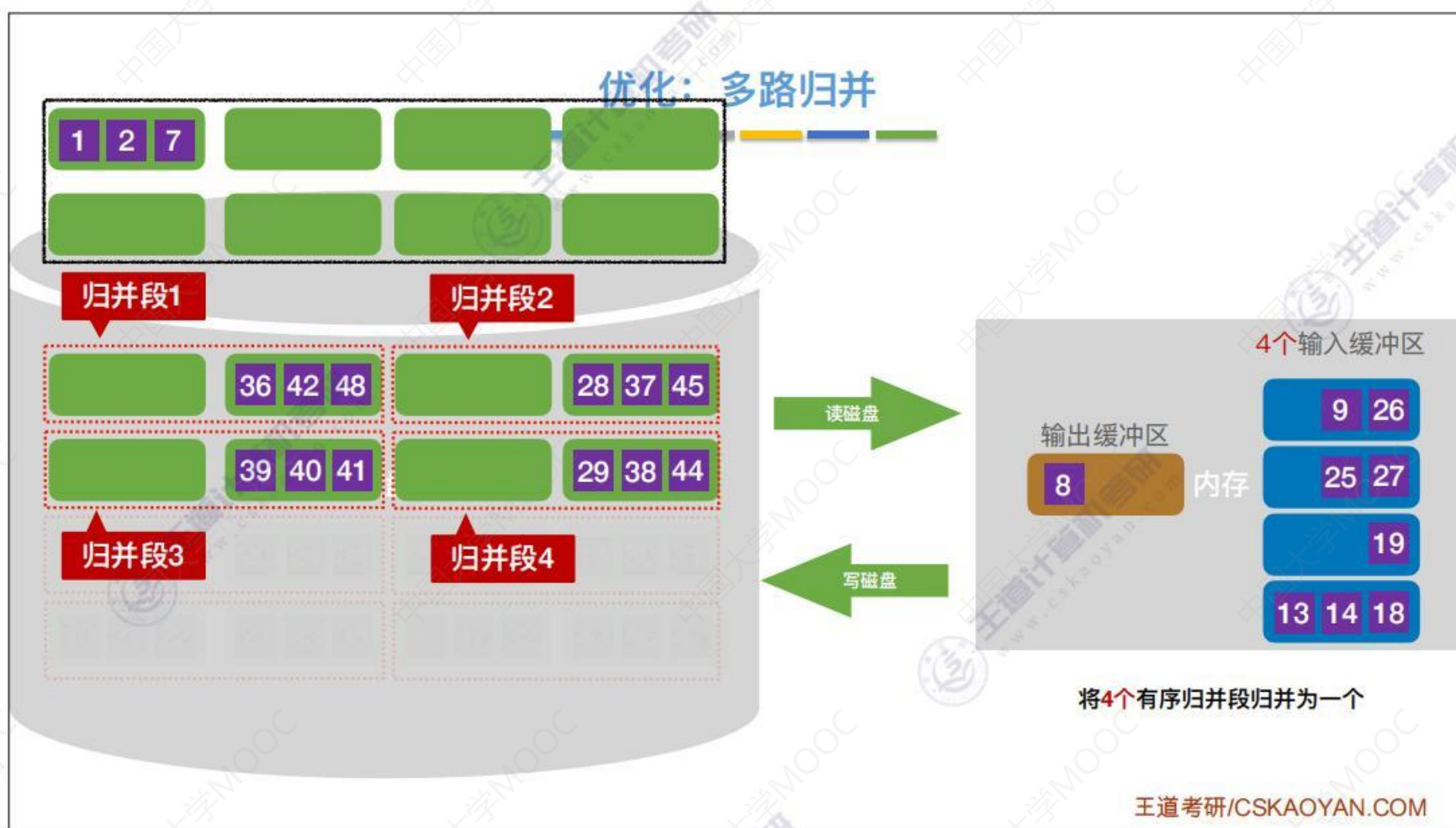


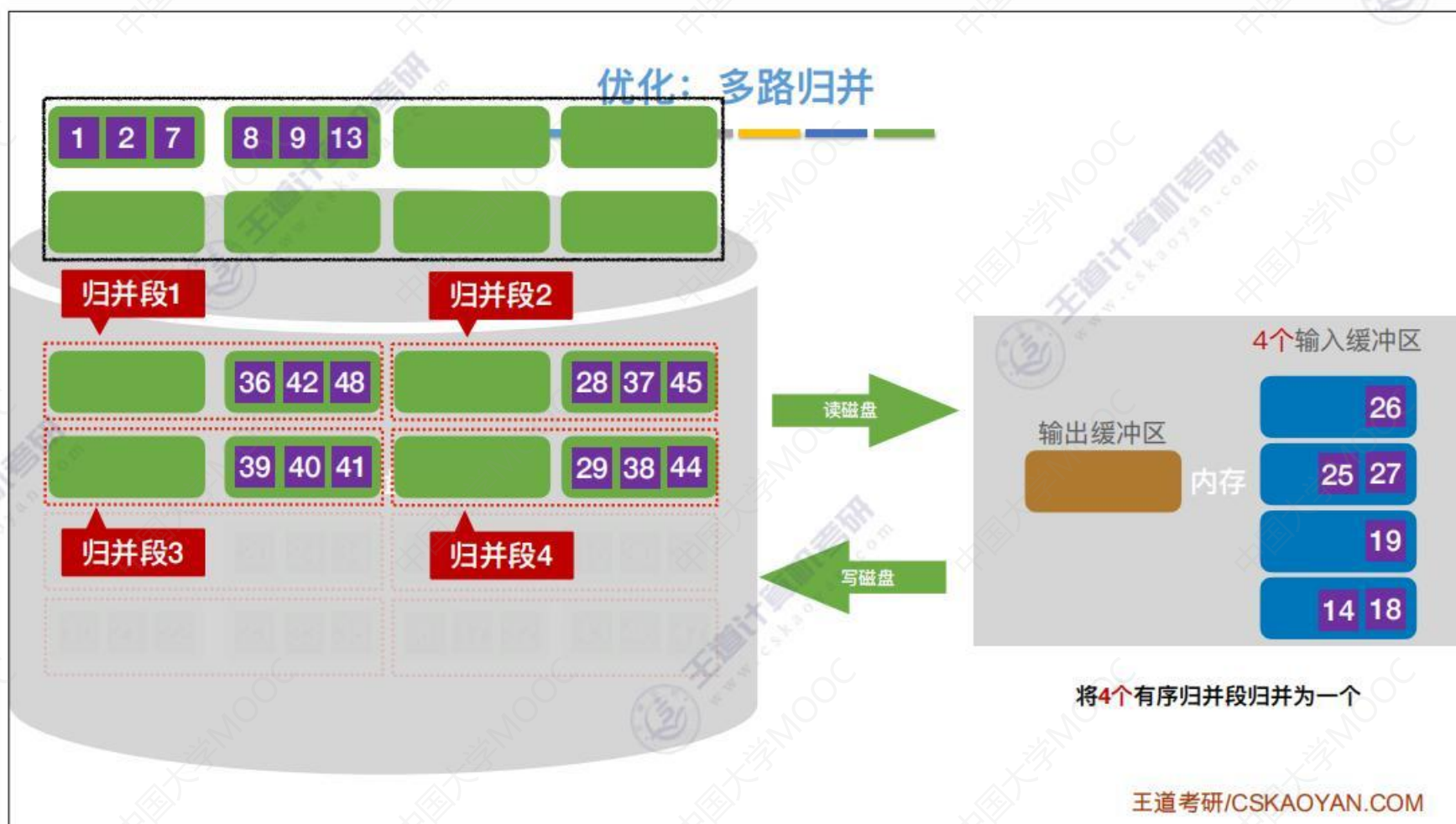
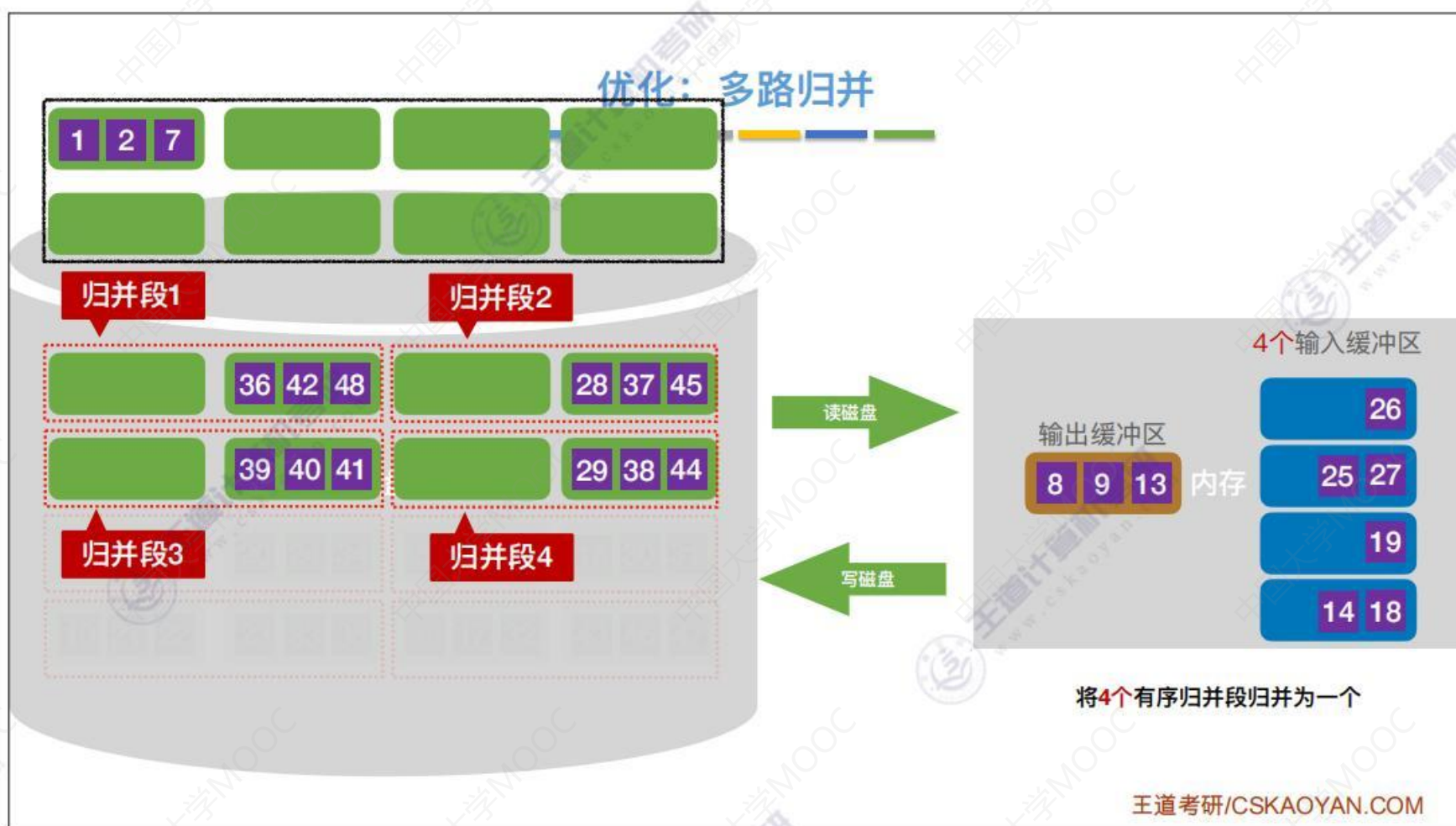
王道考研/CSKAOYAN.COM



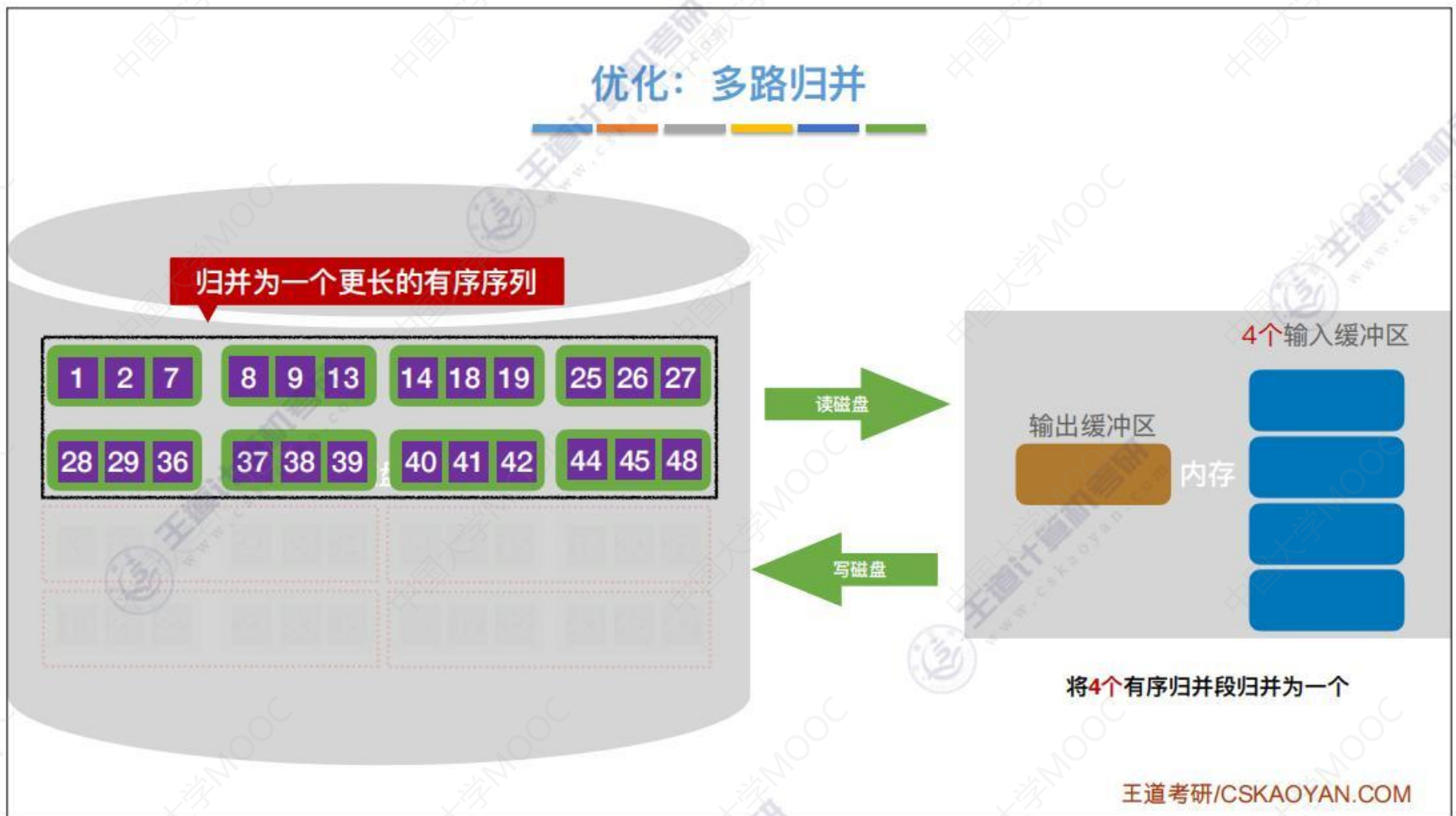




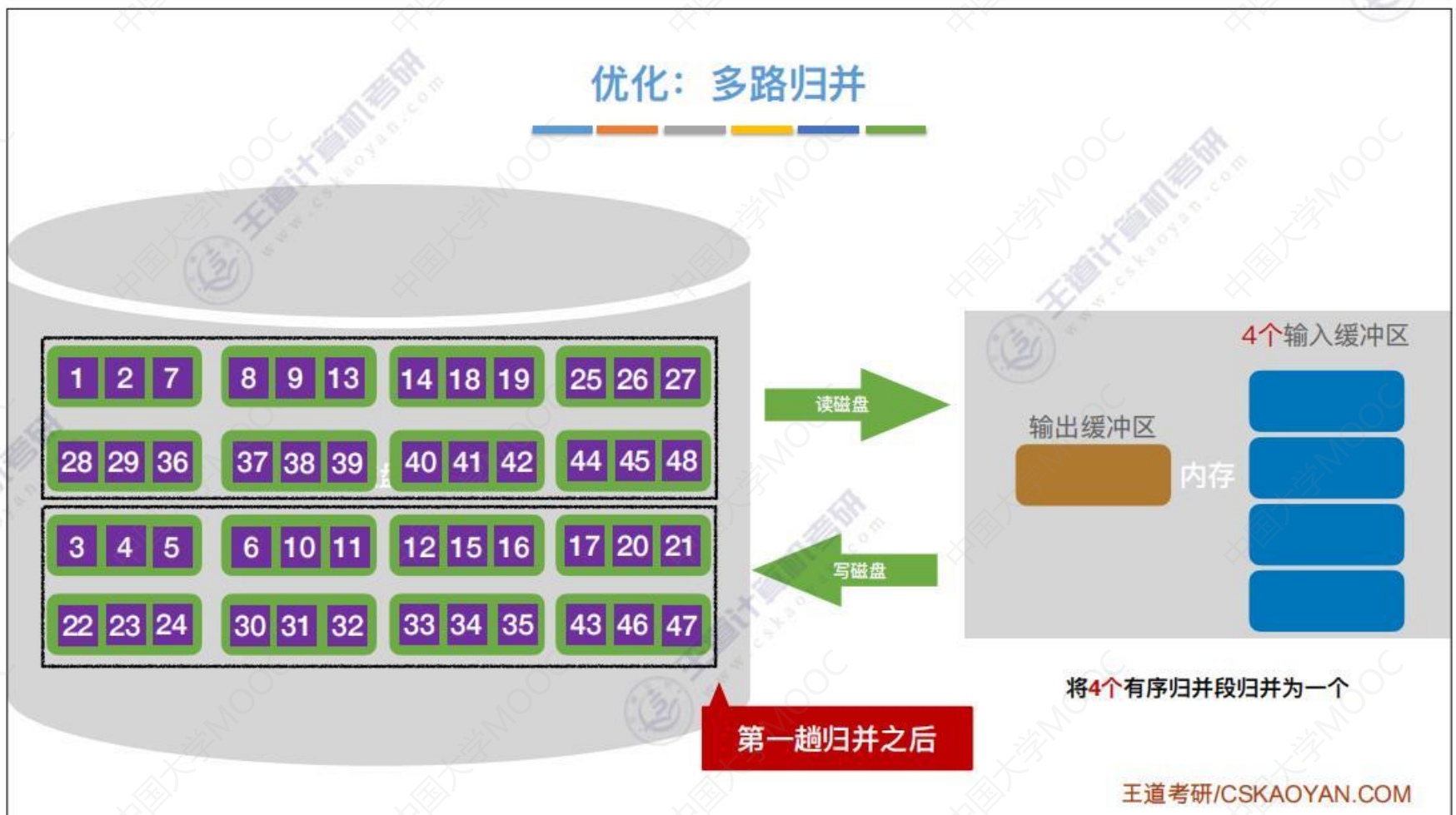




### 优化：多路归并

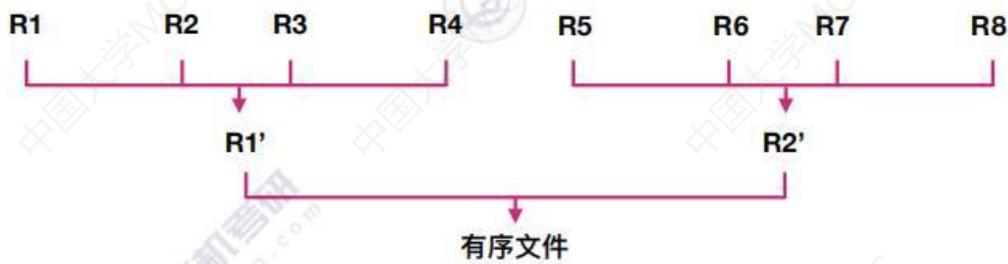


### 优化：多路归并



## 优化：多路归并

外部排序时间开销=读写外存的时间+内部排序所需时间+内部归并所需时间



采用4路归并，只需进行两趟归并即可  
读、写磁盘次数 =  $32 + 32 * 2 = 96$  次

**重要结论：**采用多路归并可以减少归并趟数，从而减少磁盘I/O(读写)次数

对  $r$  个初始归并段，做  $k$  路归并，则归并树可用  $k$  叉树表示

若树高为  $h$ ，则归并趟数 =  $h - 1 = \lceil \log_k r \rceil$

推导： $k$  叉树第  $h$  层最多有  $k^{h-1}$  个结点

则  $r \leq k^{h-1}$ ， $(h-1)_{\text{最小}} = \lceil \log_k r \rceil$

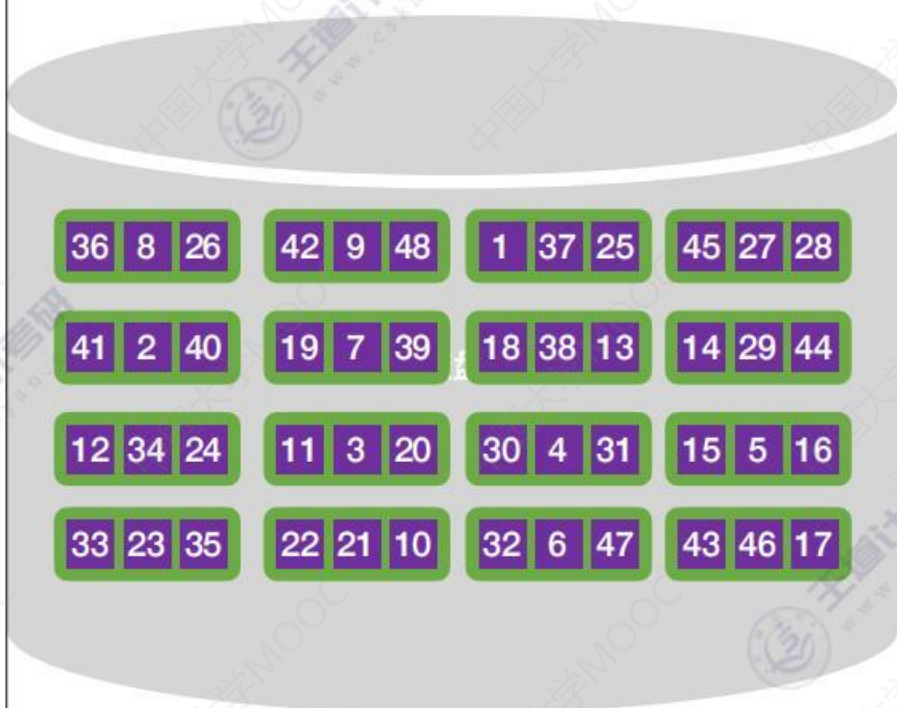
**$k$  越大， $r$  越小，归并趟数越少，读写磁盘次数越少**

多路归并带来的负面影响：

- ①  $k$  路归并时，需要开辟  $k$  个输入缓冲区，内存开销增加。
- ② 每挑选一个关键字需要对比关键字  $(k-1)$  次，内部归并所需时间增加

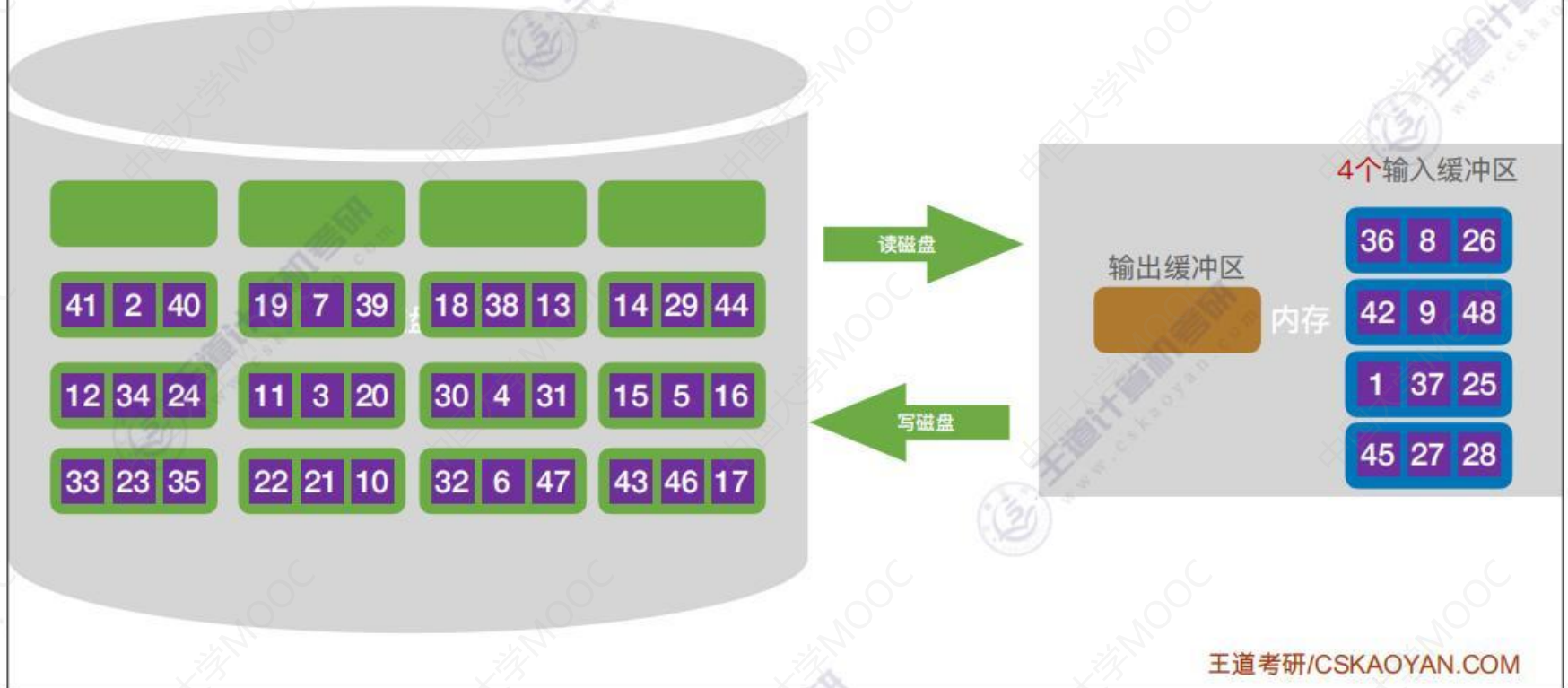
王道考研/CSKAOYAN.COM

## 优化：减少初始归并段数量

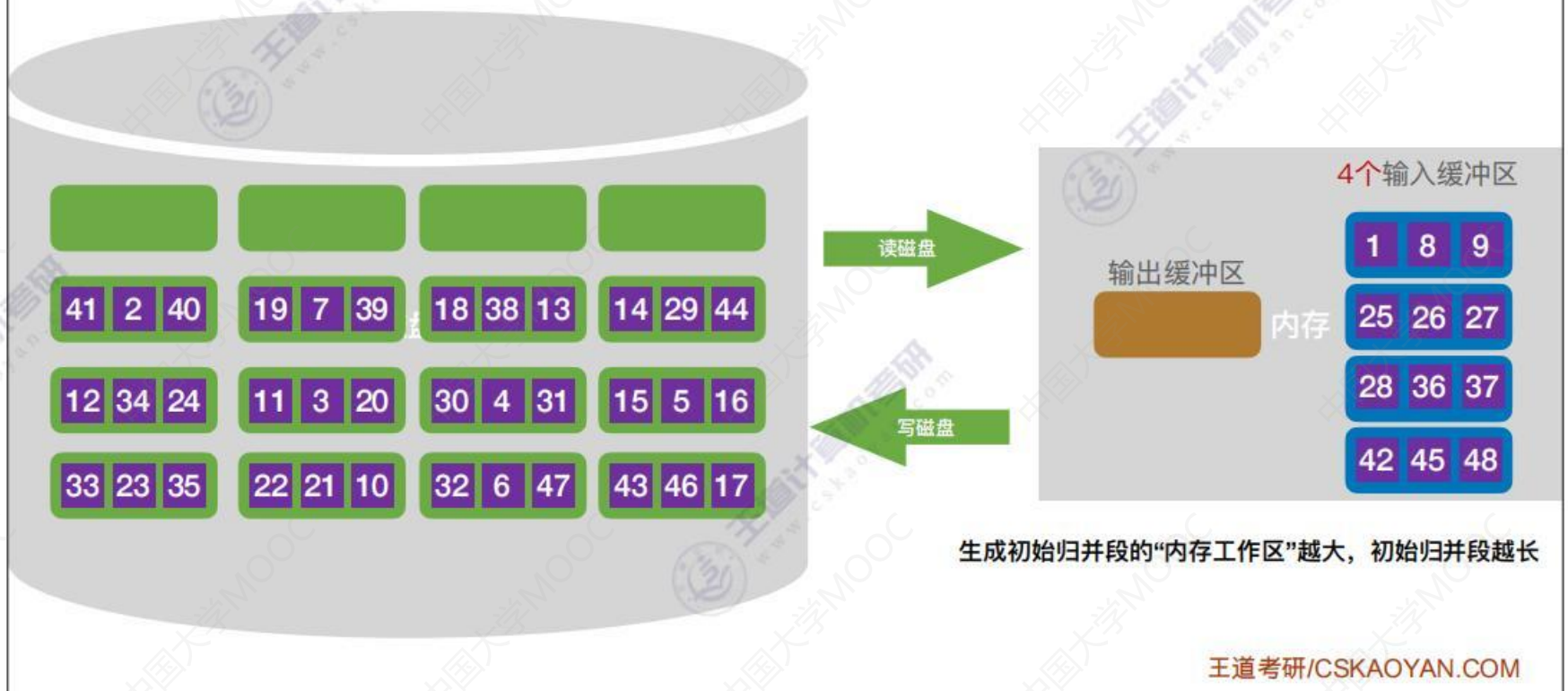


王道考研/CSKAOYAN.COM

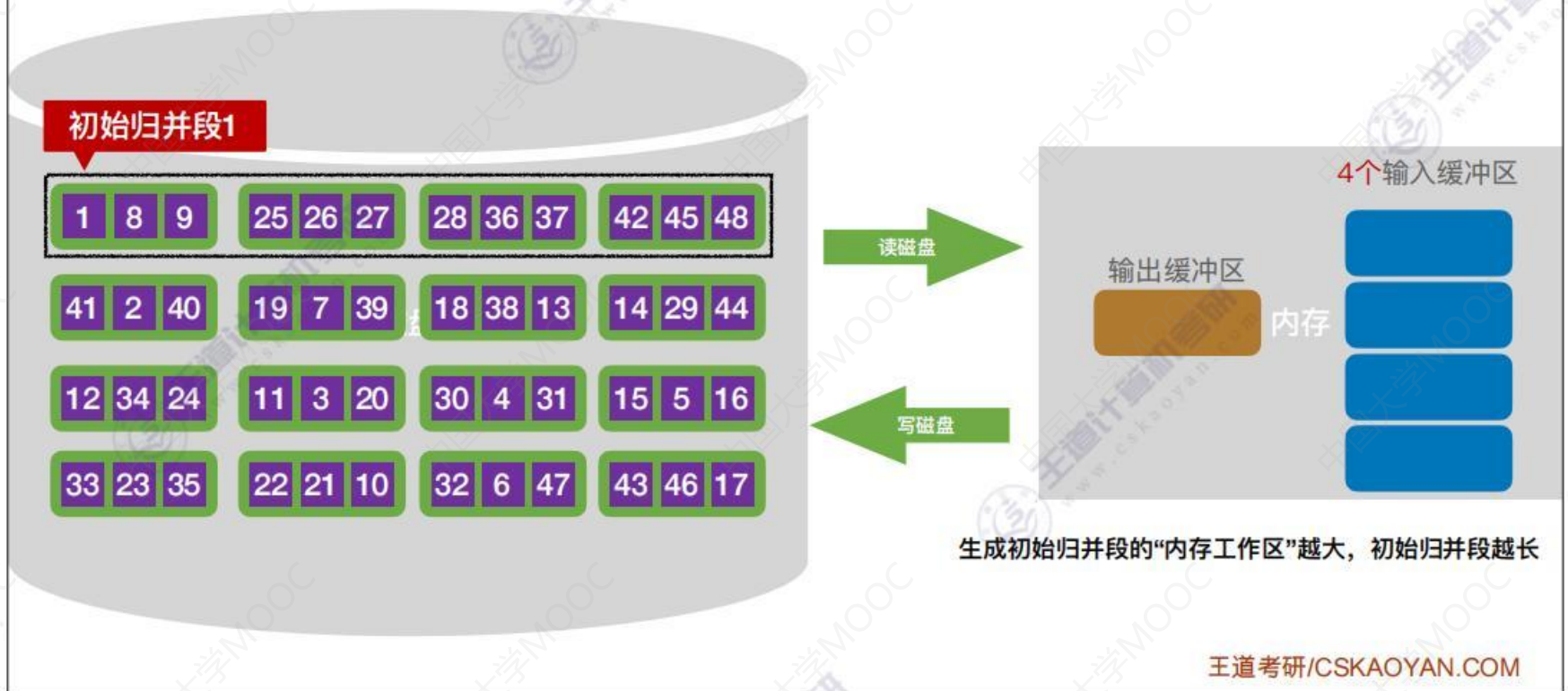
### 优化：减少初始归并段数量



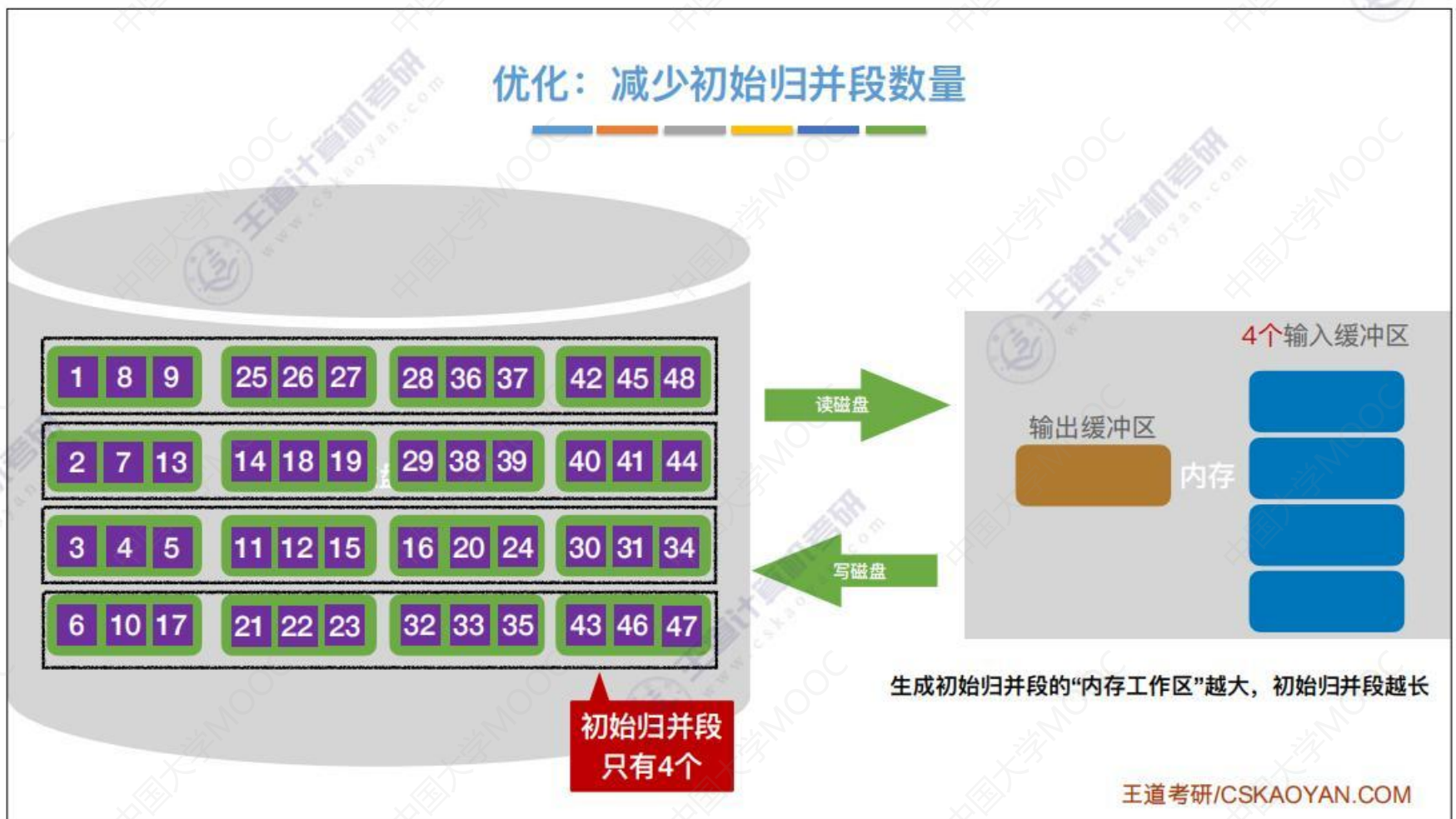
### 优化：减少初始归并段数量



### 优化：减少初始归并段数量



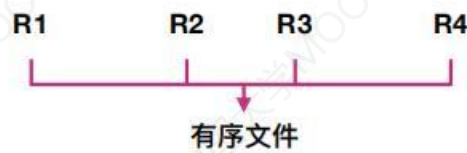
### 优化：减少初始归并段数量



## 优化：减少初始归并段数量

对  $r$  个初始归并段，做  $k$  路归并，则归并树可用  $k$  叉树表示  
若树高为  $h$ ，则归并趟数  $= h-1 = \lceil \log_k r \rceil$

$k$  越大， $r$  越小，归并趟数越少，读写磁盘次数越少

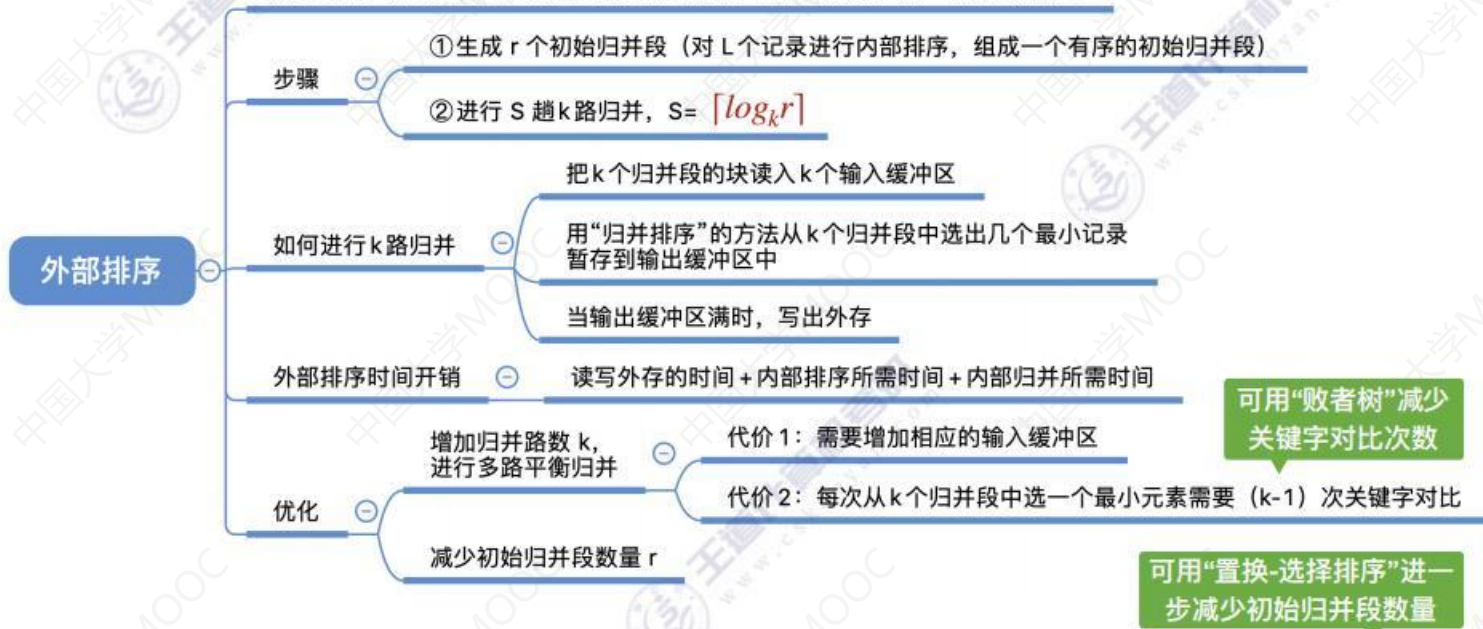


结论：若能增加初始归并段的长度，则可减少初始归并段数量  $r$

王道考研/CSKAOYAN.COM

## 知识回顾与重要考点

若要进行  $k$  路归并排序，则需要在内存中分配  $k$  个输入缓冲区和 1 个输出缓冲区



注：按照本节介绍的方法生成的初始归并段，若共  $N$  个记录，内存工作区可以容纳  $L$  个记录，则初始归并段数量  $r = \lceil N/L \rceil$

王道考研/CSKAOYAN.COM

## 纠正：什么是多路平衡归并？

误：对  $r$  个初始归并段，做  $k$  路平衡归并，归并树可用严格  $k$  叉树（即只有度为  $k$  与度为  $0$  的结点的  $k$  叉树）来表示。



知道错哪了不？

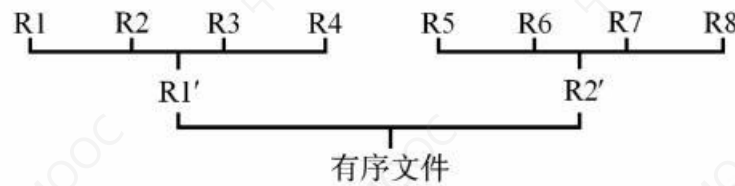


图 7.8 4 路平衡归并的排序过程

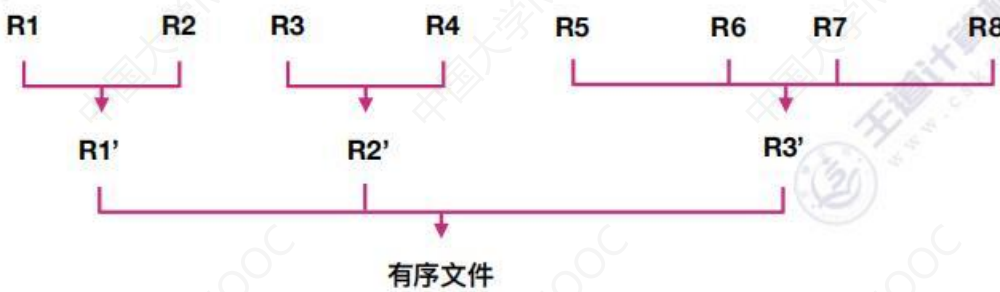
**k 路平衡归并：** ① 最多只能有  $k$  个段归并为一个；

② 每一趟归并中，若有  $m$  个归并段参与归并，则经过这一趟处理得到  $\lceil m/k \rceil$  个新的归并段

王道考研/CSKAOYAN.COM

## 纠正：什么是多路平衡归并？

8 个归并段经过一  
趟处理后得到 3 个  
新的归并段



这个例子是不是 4 路归并排序？—— 是！

这个例子是不是 4 路平衡归并排序？—— 不是！

$$\lceil 8/4 \rceil = 2 \neq 3$$

**k 路平衡归并：** ① 最多只能有  $k$  个段归并为一个；

② 每一趟归并中，若有  $m$  个归并段参与归并，则经过这一趟处理得到  $\lceil m/k \rceil$  个新的归并段

王道考研/CSKAOYAN.COM



@王道论坛



@王道计算机考研备考



@王道计算机考研



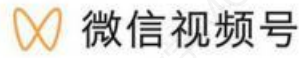
@王道咸鱼老师-计算机考研



@王道楼楼老师-计算机考研



@王道计算机考研



@王道计算机考研



@王道在线