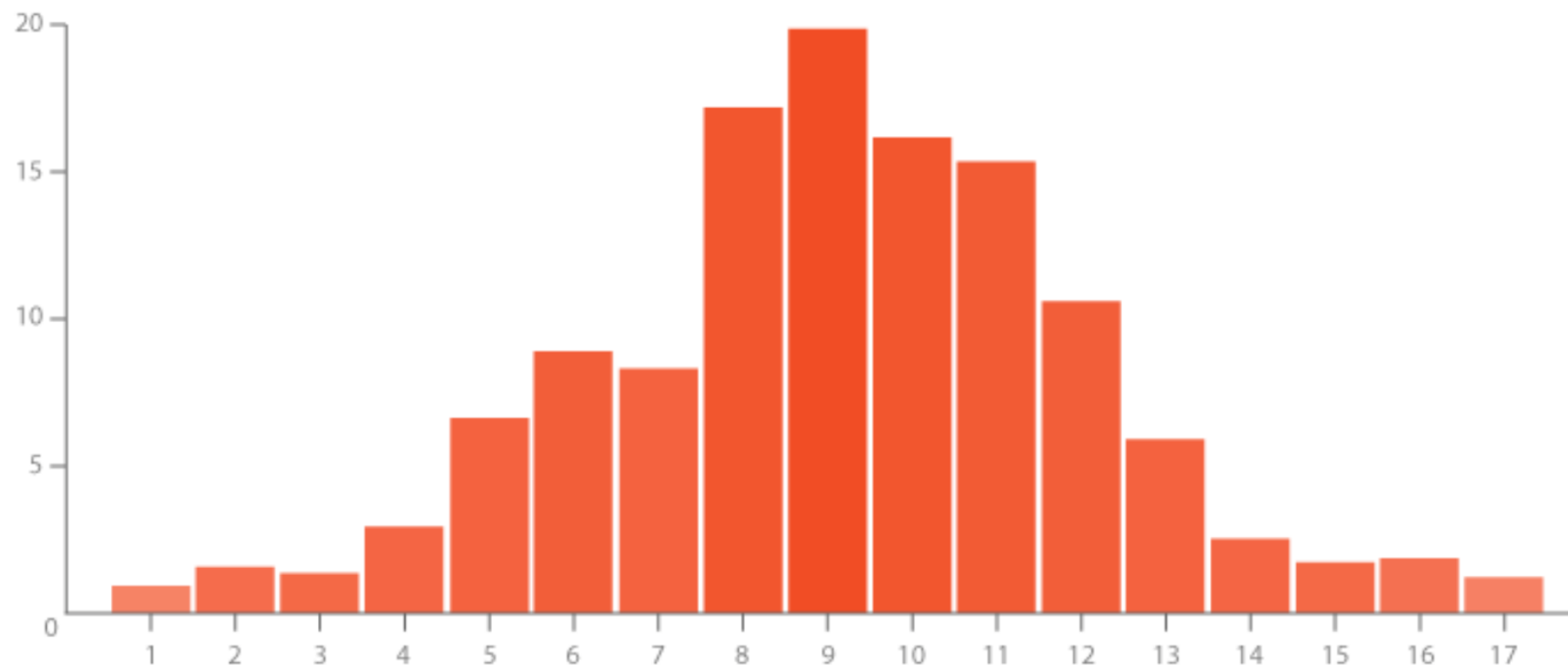


GPU并行计算与CUDA编程 第4课

- 并行化高效策略（二）
 - 1. 直方图（实例）
 - 2. 压缩与分配
 - 2.1 什么是压缩
 - 2.2 什么是分配
 - 2.3 分配的策略
 - 3. 分段扫描
 - 3.1 稀疏矩阵与向量相乘
- 4. 排序
 - 4.1 奇偶排序
 - 4.2 归并排序
 - 4.3 排序网

1. 直方图



- 方法一：直接做累加

```
__global__ void naive_histo(int *d_bins, const int *d_in, const int BIN_COUNT)
{
    int myId = threadIdx.x + blockDim.x * blockIdx.x;
    int myItem = d_in[myId];
    int myBin = myItem % BIN_COUNT;
    d_bins[myBin]++;
}
```

- 这样的做法正确吗？

- 方法二：原子相加

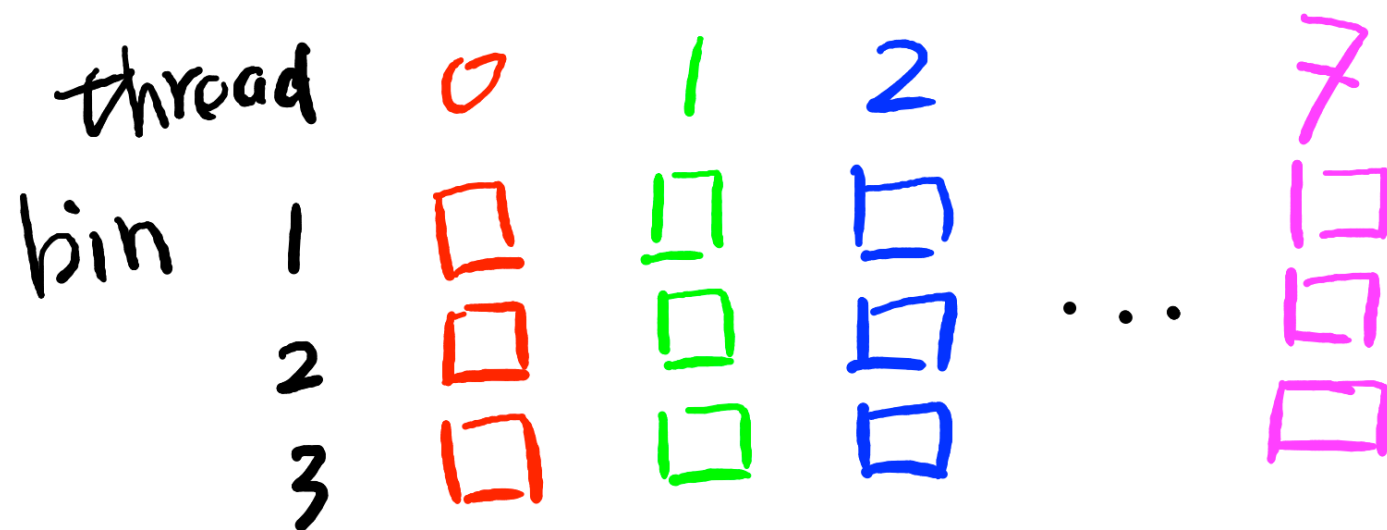
```
__global__ void simple_histo(int *d_bins, const int *d_in, const int BIN_COUNT)
{
    int myId = threadIdx.x + blockDim.x * blockIdx.x;
    int myItem = d_in[myId];
    int myBin = myItem % BIN_COUNT;
    atomicAdd(&(d_bins[myBin]), 1);
}
```

- 问题：并行化程度太低（分组bins越少，并行化程度越低，方法二适合用于分组bins很多的时候）

- 方法三：局部直方图
- 假设现有0-127的数字,用8个线程,3个分组（按除以3的余数分组）

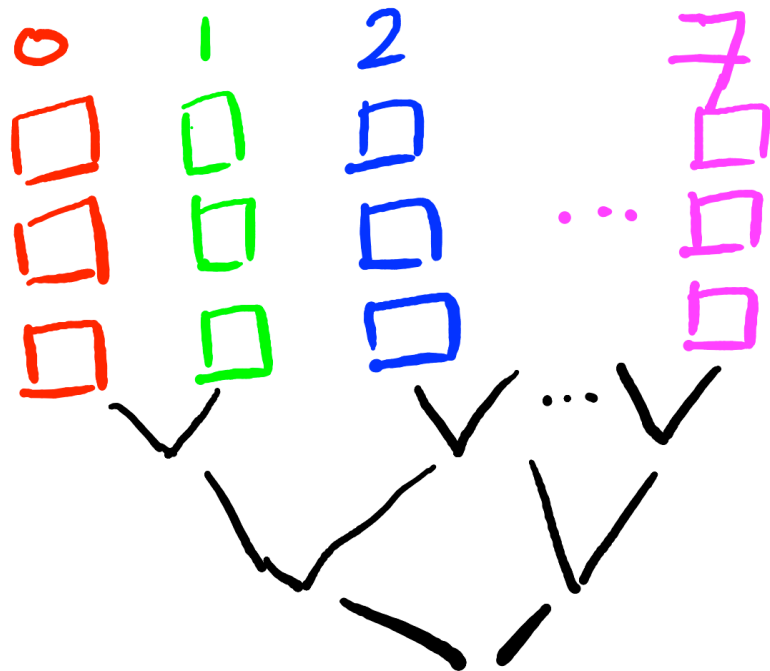
第一步：并行计算局部直方图

128 Items, 8 thread, 3 bin



- 第二步：把所有局部直方图每个分组bin使用Reduction（归约）并行累加起来行程一个总的直方图。

Reduce 8 local histo



2. 压缩与分配

- 情景：现在有一套扑克牌（52张），需要选取其中的方块的扑克牌（共13张）
- 方法一：（稀疏型）**比如说稀疏矩阵。** 方法二：（密集型）

```
if (card.isDIAMOND())  
{  
    computeCard();  
}
```

→ 52 threads

```
cc = compact(cards, isDIAMOND());  
map(cc, computeCard());
```

→ 52 threads 13 threads

疑问：这里的说法是，if和else中执行的操作是不一样的，因此执行的时间也不一样，因此可能导致有些线程闲置了，而线程块中的线程要一起运行完了才能结束，因此时间变慢。也就是防止线程发散。

问题1：if和else区别这么大么？

- 压缩 (Compact) 是什么？
- 概括的来讲，就是在输出中，对真值输入分配1，对假值输入分配0

- 压缩的步骤
- 1、 通过一个并行判断把需要进行分析的线程先通过一个判断全部筛选出来。
- 2、 把筛选好的线程做并行分析。

- 分配 (Allocate) 是什么？
- 输出项数可以动态的从每一个输入项计算出。

3. 分段扫描

(1 2 | 3 4 5 | 6 7 8)

1 3 3 7 12 6 13 21

- 例子:稀疏矩阵向量 (sparse matrix vector)
- 稀疏矩阵有很多零，我们要找一种结构把零都去掉

$$\begin{array}{l}
 \rightarrow \begin{matrix} & 0 & 1 & 2 \\ \begin{bmatrix} a & 0 & b \\ c & d & e \\ 0 & 0 & f \end{bmatrix} & \begin{bmatrix} x \\ y \\ z \end{bmatrix} \end{matrix} \\
 \rightarrow \\
 \rightarrow
 \end{array}$$

VALUE [a b c d e f]
 COLUMN [0 2 0 1 2 1]
 ROWPTR [0 2 5]

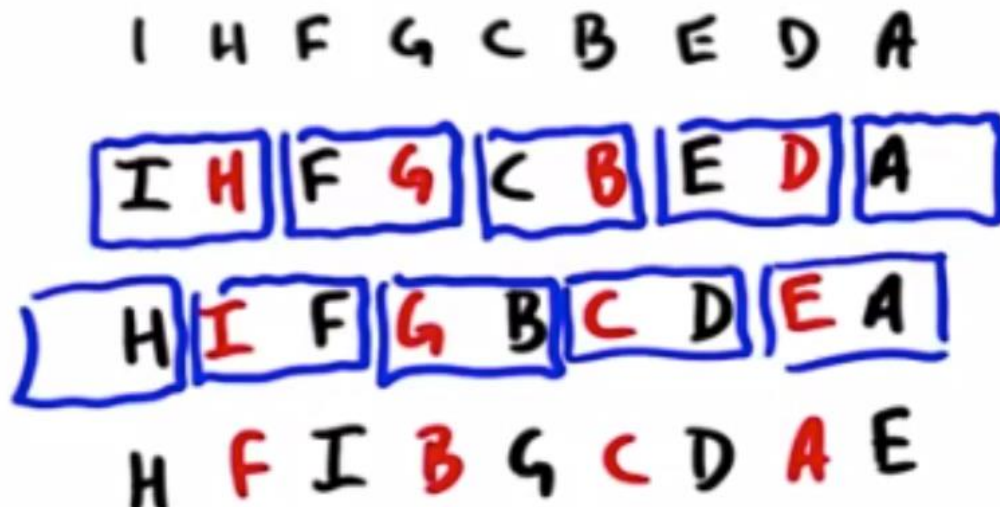
$$\begin{array}{l}
 \begin{bmatrix} a & b & | & c & d & e & | & f \end{bmatrix} \\
 \uparrow \qquad \qquad \uparrow \qquad \qquad \qquad \uparrow \\
 \begin{bmatrix} x & z & & x & y & z & & y \end{bmatrix} \\
 \\
 \begin{bmatrix} a \cdot x & b \cdot z & | & c \cdot x & d \cdot y & e \cdot z & | & f \cdot y \end{bmatrix} \\
 \underbrace{\qquad \qquad \qquad + \qquad \qquad \qquad} \qquad \underbrace{\qquad \qquad \qquad - \qquad \qquad \qquad} \qquad \underbrace{\qquad \qquad \qquad} \\
 \text{out}(0) \qquad \qquad \qquad \text{out}(1) \qquad \text{out}(2)
 \end{array}$$

4. 排序

- 难点：大部分的排序算法都是串行的
- 适合并行的排序算法：
- 首先考虑能合并内存并且线程分支较少的并行算法，尤其是能让线程保持同时忙碌的算法

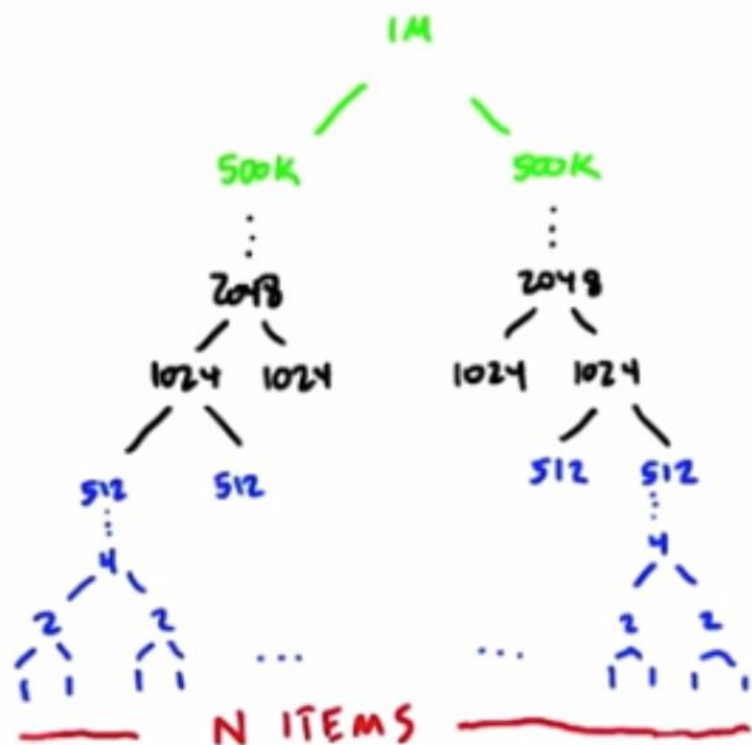
- 1. 奇偶排序（砖排序，冒泡算法）
- 通过比较数组中相邻的（奇-偶）位置数字对，如果该奇偶对是错误的顺序（第一个大于第二个），则交换。下一步重复该操作，但针对所有的（偶-奇）位置数字对。如此交替进行下去。

ODD-EVEN SORT (BRICK SORT)



- 2.归并排序
- 将两个已经排序的序列合并成一个序列的操作。

MERGE SORT



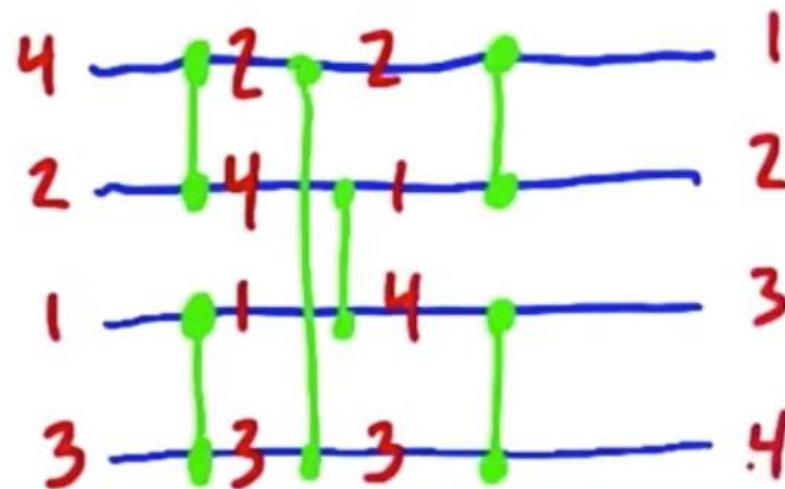
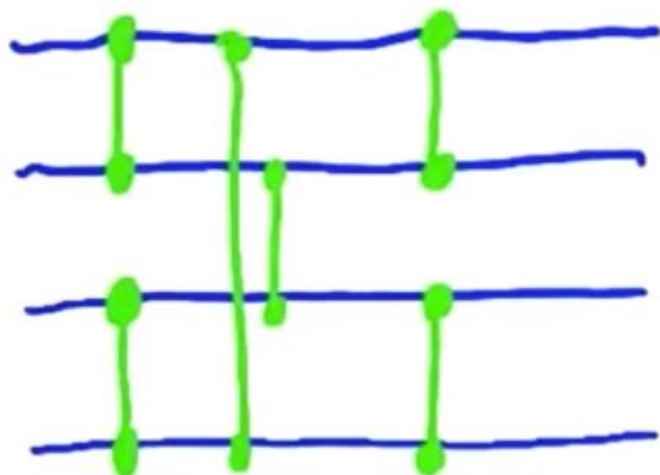
STAGE 3
ONE TASK
BIG TASK!

SPLIT TASK ACROSS SMS

STAGE 2
BUNCH OF TASKS
EACH TASK: MEDIUM
TASK PER BLOCK

STAGE 1
TONS OF TASKS
EACH TASK: SMALL
TASK PER THREAD

3. 排序网



本周作业

🔗 自行编写代码实现并行计算直方图的“方法三”，上传代码和截图效果。

【声明】 本视频和幻灯片为炼数成金网络课程的教学资料，
所有资料只能在课程内使用，不得在课程以外范围散播，
违者将可能被追究法律和经济责任。

课程详情访问炼数成金培训网站

<http://edu.dataguru.cn>

- Dataguru（炼数成金）是专业数据分析网站，提供教育，媒体，内容，社区，出版，数据分析业务等服务。我们的课程采用新兴的互联网教育形式，独创地发展了逆向收费式网络培训课程模式。既继承传统教育重学习氛围，重竞争压力的特点，同时又发挥互联网的威力打破时空限制，把天南地北志同道合的朋友组织在一起交流学习，使到原先孤立的学习个体组合成有组织的探索力量。并且把原先动辄成千上万的学习成本，直线下降至百元范围，造福大众。我们的目标是：低成本传播高价值知识，构架中国第一的网上知识流转阵地。
- 关于逆向收费式网络的详情，请看我们的培训网站 <http://edu.dataguru.cn>

Thanks

FAQ时间