Now, this is just a simulation of what the blocks will look like once they've assembled
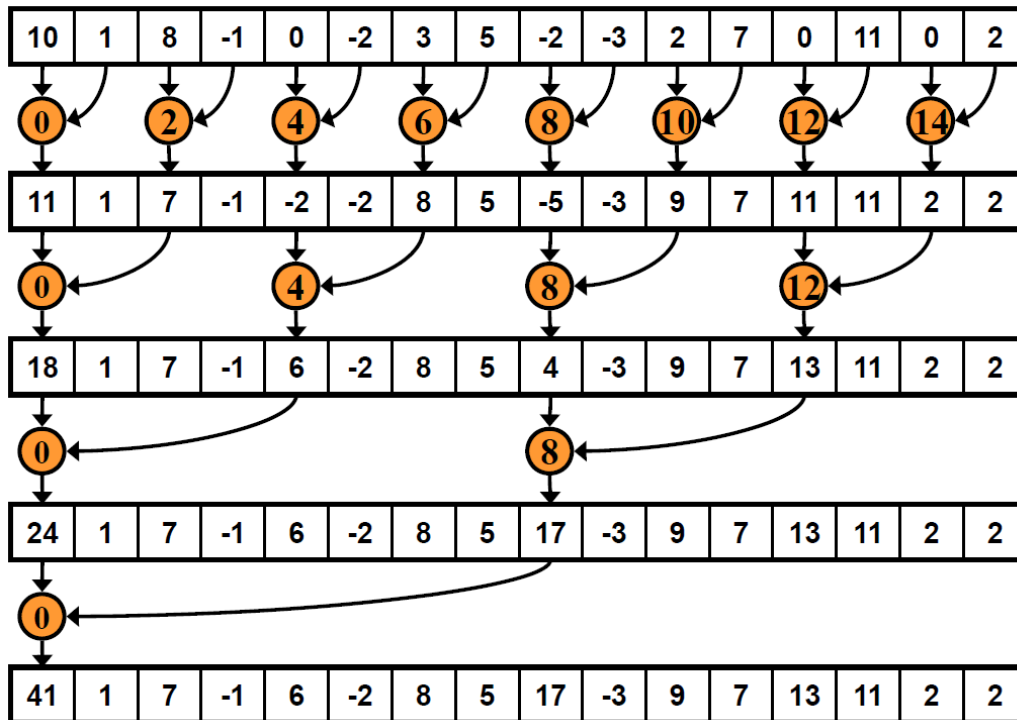
- Reduction – Stride Pattern – Shared Memory and Branch Divergence

# Reduction - Tesla P100;compute v6.0;

n: 1<<20

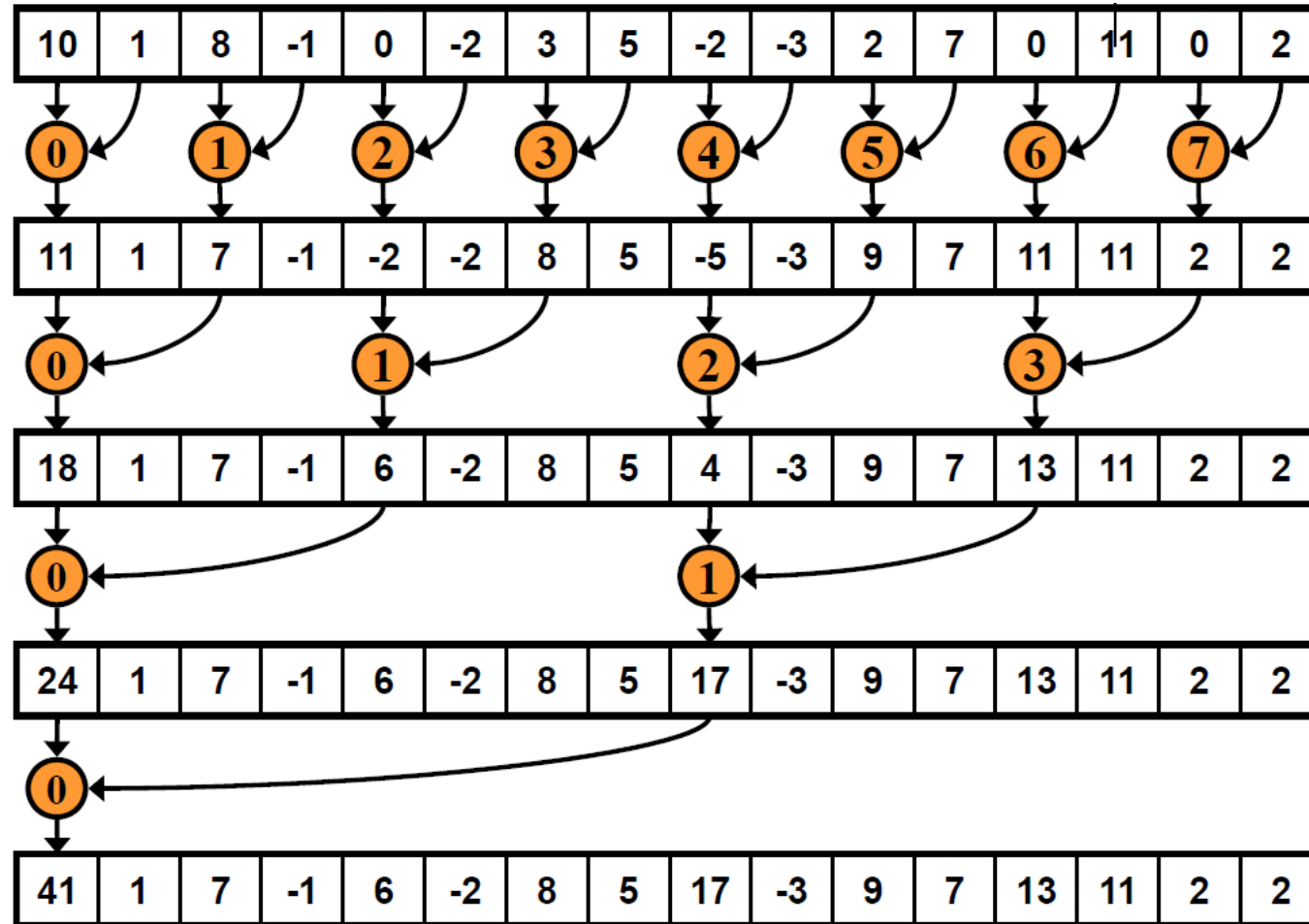| Version | Time (ms) |
|---|---|
| serial | 3.27400 |
| global reduce stride – naïve | 0.16450 |
| shared stride reduce | 0.15835 |

20.7X



Unusual we should have had better performance!

How to reorganize workload assignment to avoid divergence?
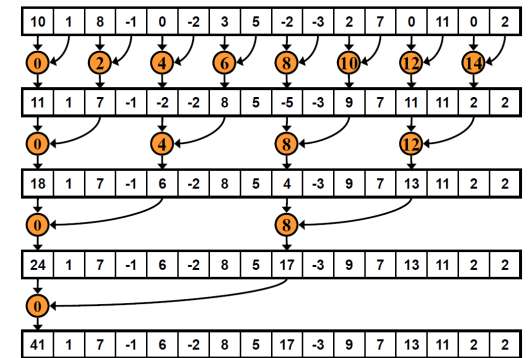
# Kernel: Shared Memory – Stride Pattern

# Kernel: Shared Memory – Stride Pattern

```
__global__ void shared_reduce_stride(float* d_out, float* d_in){
  extern __shared__ float sdata[];
  // shared_reduce<<<blocks,threads,threads*sizeof(float)>>>
  int myId = threadIdx.x + blockDim.x * blockIdx.x;
  int tid  = threadIdx.x;
  // load shared mem from global mem
  sdata[tid] = d_in[myId];
  // make sure entire block is loaded!
  // do reduction in shared memory
  for(int stride = 1; stride < blockDim.x; stride *= 2)  {
     __syncthreads();
     if(myId % (2*stride) == 0) {
       sdata[tid] += sdata[tid+stride]; }
     }
// thread 0 writes result for this block back to global mem
 if (tid == 0) {
       d_out[blockIdx.x] = sdata[tid]; }
}
```
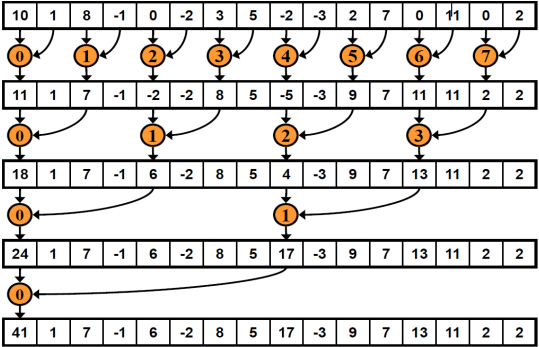
# Kernel: Shared Memory – No Diverge – Stride Pattern

```
__global__ void shared_reduce_stride_nodiverge (float* d_out,
float* d_in){
  extern __shared__ float sdata[];
  // shared_reduce<<<blocks,threads,threads*sizeof(float)>>>
  int myId = threadIdx.x + blockDim.x * blockIdx.x;
  int tid  = threadIdx.x;
  // load shared mem from global mem
  sdata[tid] = d_in[myId];
  // make sure entire block is loaded!
  // do reduction in shared memory
  for(int stride = 1; stride < blockDim.x; stride *= 2)  {
      __syncthreads();



  }
  // thread 0 writes result for this block back to global mem
  if (tid == 0) {
        d_out[blockIdx.x] = sdata[tid]; }
}
```
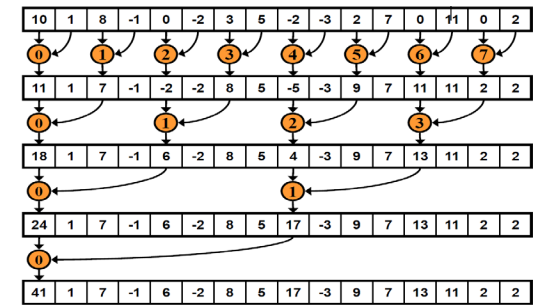
# Kernel: Shared Memory – No Diverge – Stride Pattern

```
__global__ void shared_reduce_stride_nodiverge (float* d_out,
float* d_in){
    extern __shared__ float sdata[];
    // shared_reduce<<<blocks,threads,threads*sizeof(float)>>>
    int myId = threadIdx.x + blockDim.x * blockIdx.x;
    int tid  = threadIdx.x;
    // load shared mem from global mem
    sdata[tid] = d_in[myId];
    // make sure entire block is loaded!
    // do reduction in shared memory
    for(int stride = 1; stride < blockDim.x; stride *= 2)  {
        __syncthreads();
        int index = 2*stride*tid;
        if( index < blockDim.x ) {
            sdata[index] += sdata[stride+index]; }
    }
// thread 0 writes result for this block back to global mem
 if (tid == 0) {
        d_out[blockIdx.x] = sdata[tid]; }
}
```
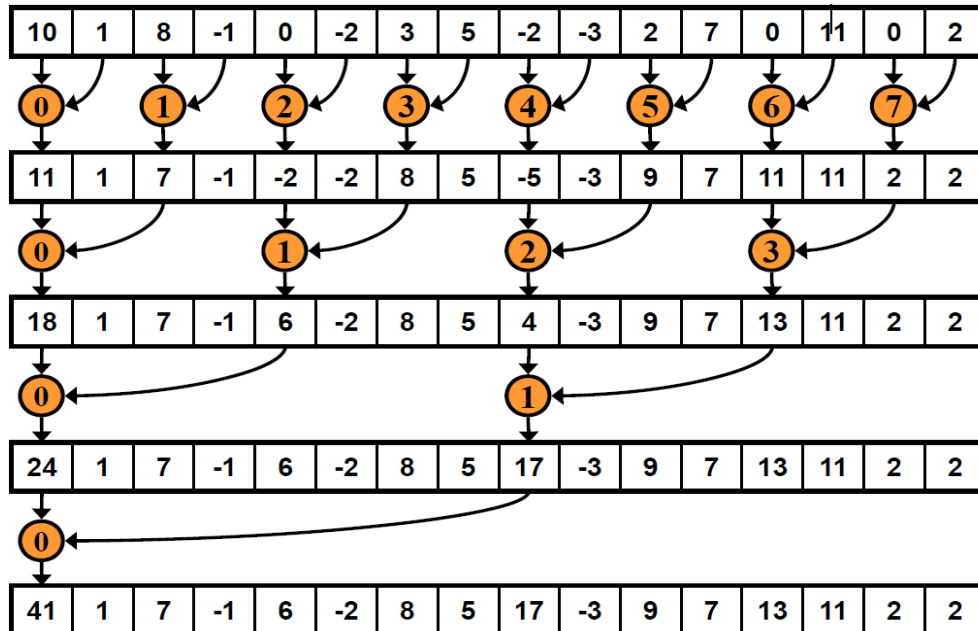
Bonus: No more expensive % operator

# Reduction - Tesla P100;compute v6.0;

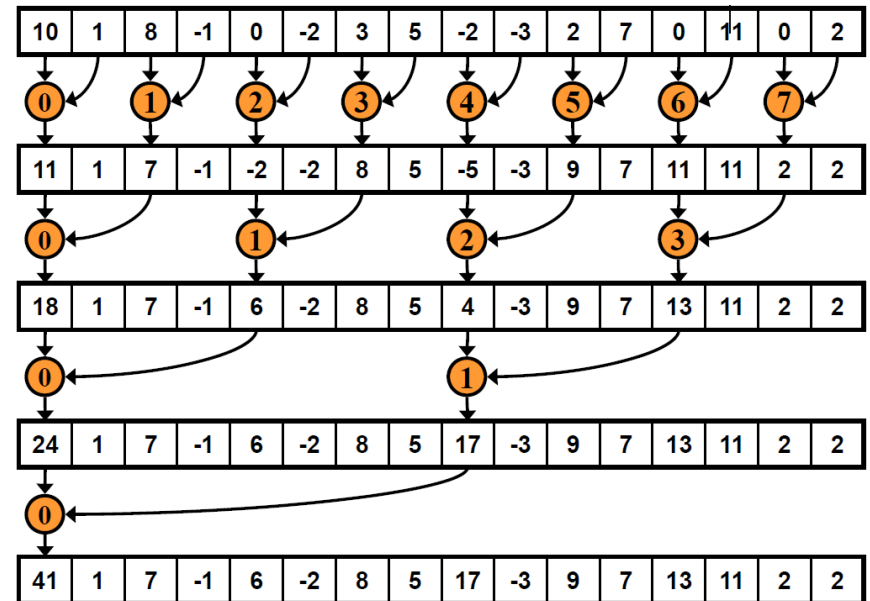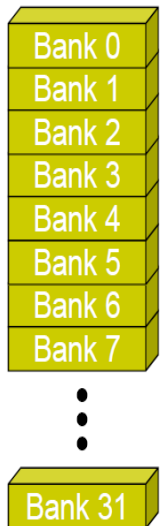| Version | Time (ms) |
|---|---|
| serial | 3.27400 |
| global reduce stride – naïve | 0.16450 |
| shared stride reduce | 0.15835 |
| shared_reduce_stride_nodiverge | 0.09081 |

n: 1<<20

36.1X



1.7X!

# Observation on Stride Pattern-No divergence

- **Divergence free**
- **New problem:**
  - **Shared memory bank conflicts**
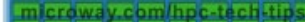  - Will come back to reduction, first bank conflicts!

# Shared Memory and Bank Conflicts

- Many threads access memory at the same time
  - To service more than one thread, memory is divided into independent banks
  - This layout essential to achieve high bandwidth
- Each SM has Shared Memory organized in 32 Memory banks

Bank 0
Bank 1
Bank 2
Bank 3
Bank 4
Bank 5
Bank 6
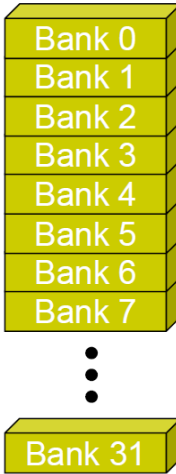Bank 7

Bank 31

# Shared Memory Architecture

- The 32 banks of the Shared Memory are organized like benches in a movie theater
  - You have multiple rows of benches
  - Each row has 32 benches which are separated and grouped in long columns
  - In each bench you can "seat" a family of four bytes (32bits total)
  - Note that a bank represents a column of benches in the movie theater, which is perpendicular to the screen

# Organization of Memory Banks



32 threads in a warp attempt to access shared memory simultaneously

Bank conflict: the scenario where two different threads access *different* words in the same bank

- -No "bank conflicts" between threads of different warps

# Bank Conflict Example

- Bank = (address of offset) % 32

- Example:

- 1D shared mem array, myShMemVar, of 1024 floats
  - myShMemVar[4]: accesses bank #4 (the fifth one – first row)
  - myShMemVar[31]: accesses bank #31 (the last one – first row)
  - **myShMemVar[50]:** access bank #18 (the 19th one – second row)
  - myShMemVar[128]: access bank #0 (the first one – fifth row)
  - **myShMemVar[178]:** access bank #18 (the 19th one – sixth row)

- NOTE: If, for instance, the third thread in a warp accesses myShMemVar[50] and the eight thread in the warp access myShMemVar[178], then you have a two-way bank conflict and the two transactions get serialized

Bank 0
Bank 1
Bank 2
Bank 3
Bank 4
Bank 5
Bank 6
Bank 7

Bank 31

# Bank Conflict Examples(1): 4 Byte Words
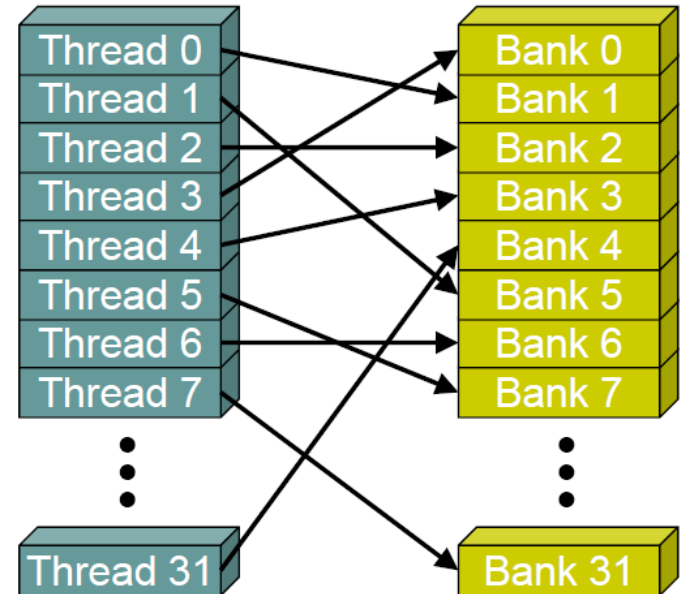
Bank = (address of offset) % 32

- ## No Bank Conflicts
  - Linear addressing stride == 1

- ## No Bank Conflicts
  - Random 1:1 Permutation

# Bank Conflict Examples(2): 4 Byte Words

## Bank = (address of offset) % 32

- sdata[threadIdx.x*2]++;

- sdata[threadIdx.x*8]++;

**Which stride amount results with bank conflicts for all threads in a warp?**

# Bank Conflict Examples(3): 4 Byte Words

Which accesses have bank conflicts?

# Bank Conflicts – Demo

```
__global__ void mykernel1(unsigned long long* time){
__shared__ float shared[1024];
// clock returns clock ticks  unsigned long long
startTime = clock();
//all threads accessing the same location (broadcast)
shared[0]++;
unsigned long long finishTime = clock();  *time =
(finishTime-startTime);
}
Replace shared[0]++ with
Kernel 2: shared[threadIdx.x]++;
Kernel 3: shared[threadIdx.x*4]++;
Kernel 4: shared[threadIdx.x*8]++;
Kernel 5: shared[threadIdx.x*32]++;
```

Refer to D2L→Demo→8.BankConflict

# Kernel: Shared  Memory – Stride Pattern

## Next: How to eliminate bank conflicts?