



PARALLEL PROCESSING

UNIT 5

Dr. Ahmed Sallam

www.sallam.cf > (sallamah.weebly.com)

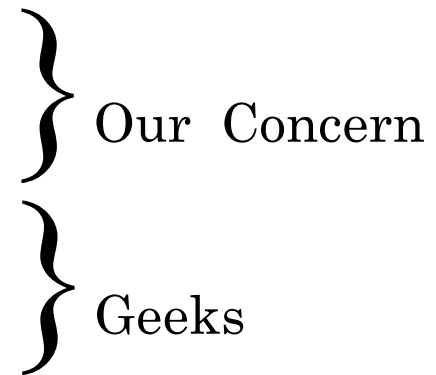
QUIZ

Principals of Efficient GPU Programming

- Decrease arithmetic intensity
- Decrease time spent on memory operations
- Coalesce global memory accesses
- Move all data to shared memory
- Do fewer memory operations per thread
- Avoid thread Divergence

LEVELS OF OPTIMIZATION

1. Picking good Algorithms
2. Basic principles for efficiency
3. Arch-specific detailed optimization
4. Instructions micro optimization



CPU: 80%

GPU: 30%

○ Quiz

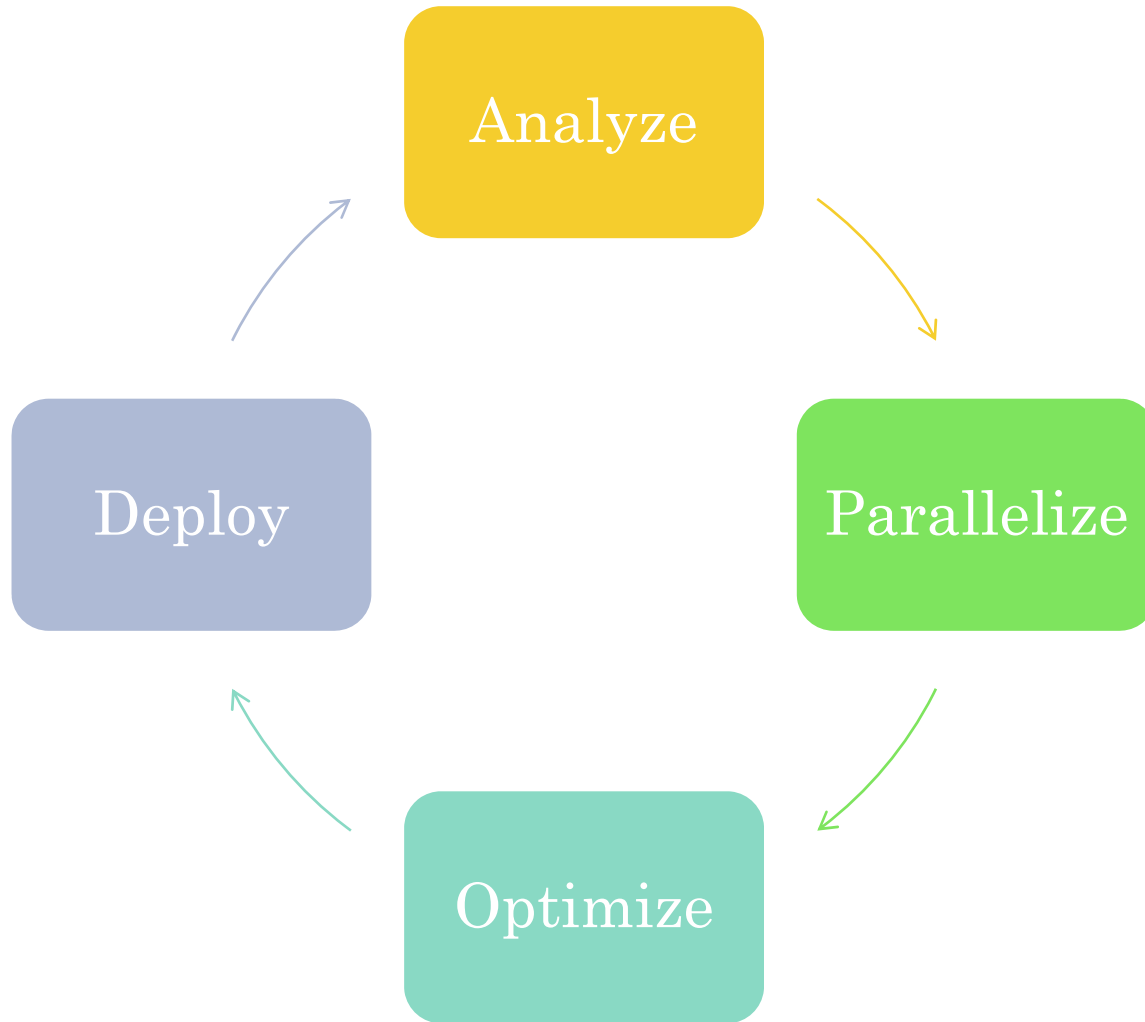
- 3 Vector registers – SSE, AVX
- 1 Use Odd-Even $O(n^2)$, Merge Sort $O(n \log n)$, Or **Heap sort** $(n \log n)$
- 2 Write cache aware code (Iterate 2D array by column)
- 3 Block for L1 cache
- 4 Perform $x^{-1/2}$ by: `0x5f3759df - (x >> 1)` (Quake3 game)



APOD SYSTEMATIC OPTIMIZATION

4

APOD



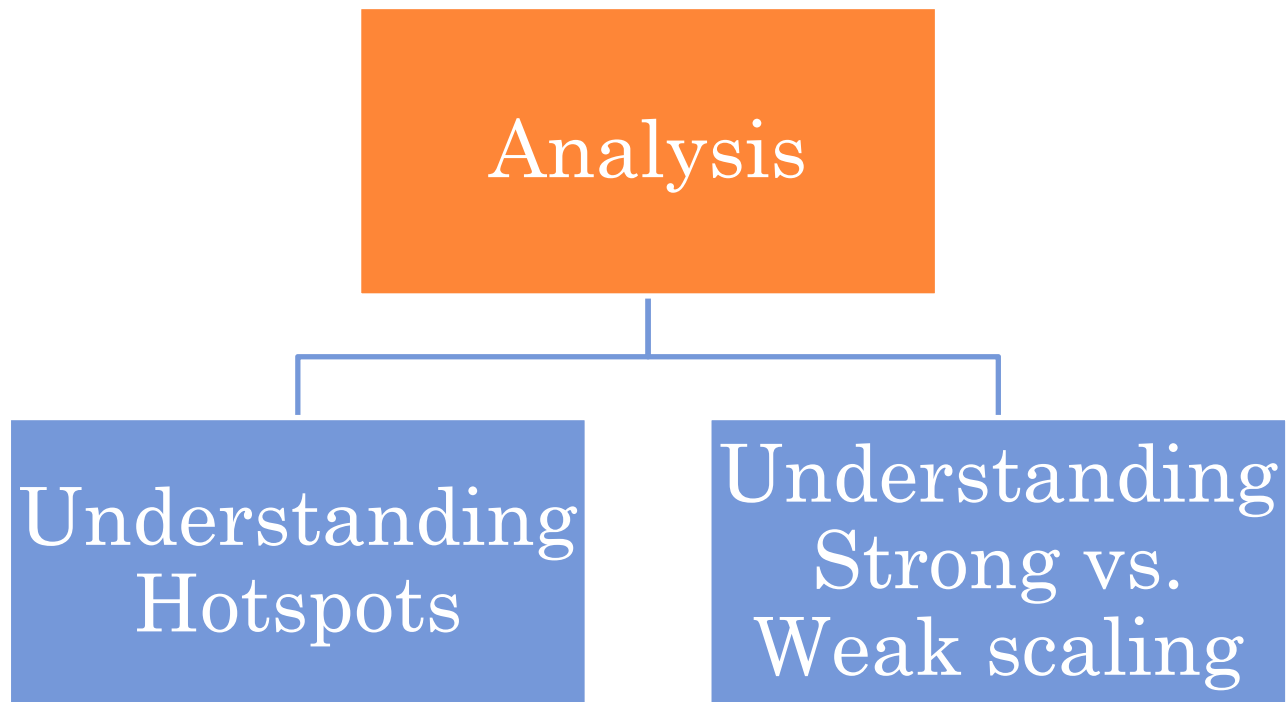
APOD (CONT.)

- Analyze: Profile whole application
 - Where can it benefit?
 - By how much?
- Parallelize:
 - Pick an Approach
 - Complete algorithms: find Libraries or write by yourself (CUDA or OpenCl)
 - Small code: Directives (OpenMP for CPU), (OpenACC for GPU)
 - Pick an algorithm
- Optimization: Profile-driven Optimization
 - (Measure don't just guess)
- Deploy: Optimize to a limit
 - don't waste your time by more optimization.

The background features a black field with vertical stripes of varying widths and textures. On the left, there are stripes with a fine grid pattern and a solid orange stripe. On the right, there is a solid orange stripe. Several orange circles of different sizes are scattered across the page, with one large circle on the left and another on the right. The word 'ANALYSIS' is written in a serif font in the center-right area.

ANALYSIS

7

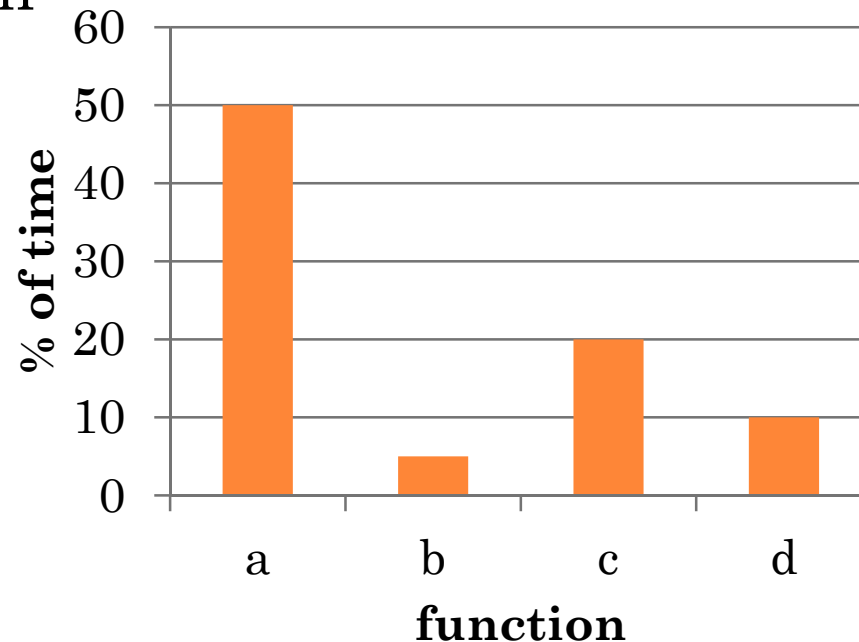


WEAK VS. STRONG SCALING

- Serial program solves problem size A in time T
 - Ex: PageRank matrix multiplication in an hour
- Weak Scaling: run a larger problem, or more small one.
- Strong Scaling: run the same size problem more faster.

UNDERSTANDING HOTSPOTS

- Don't rely on intuition
- Run a profiler
 - Gprof
 - VTune
 - Very Sleepy



- Max speed up if you parallelize (Drop to zero):
 - Speedup by parallelize function (a): = $2x$
 - Speedup by parallelize functions (a, c, d): = $5x$

UNDERSTANDING HOTSPOTS (CONT.)

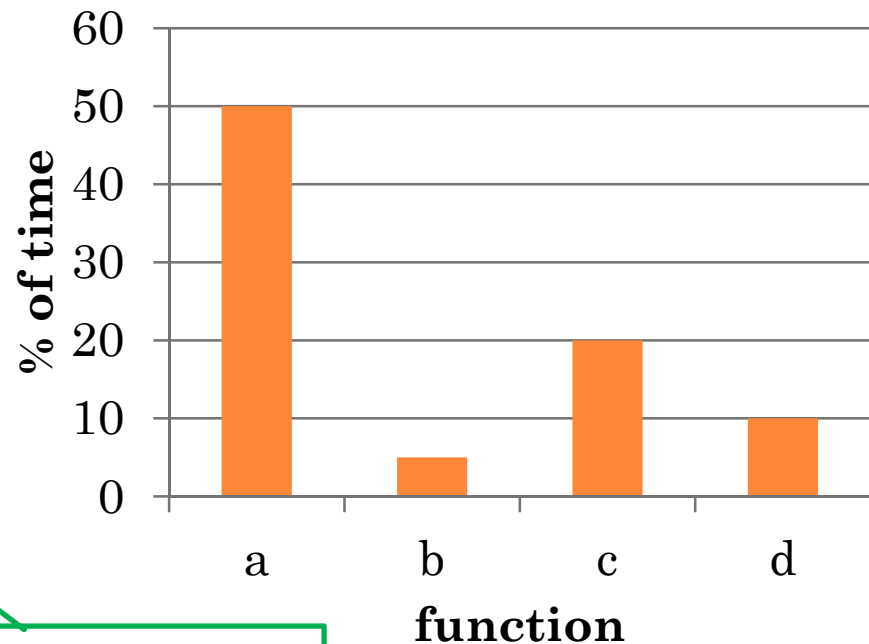
- Amdahl's Law

- Total speedup from parallelization is limited by portion of time spent doing something to be parallelized.

- $\text{max speedup} = \frac{1}{s + \frac{p}{n}}$

where s is serial portion and $p = (1 - s)$ is parallel portion, n is the number of threads

$$= \frac{1}{s} = \frac{1}{1-p}$$



We drop this part to 0 because we have enough n

- Previous example:

- $p(a): 1/(1-50\%) = 2x$
- $p(a, c, d): 1/(1-80\%) = 5x$

The background features a black field with vertical stripes of varying widths and colors, including a wide textured grey stripe on the left and a thin orange stripe on the right. Several orange circles of different sizes are scattered across the page, with one large circle on the left and another on the right.

PARALLELIZE

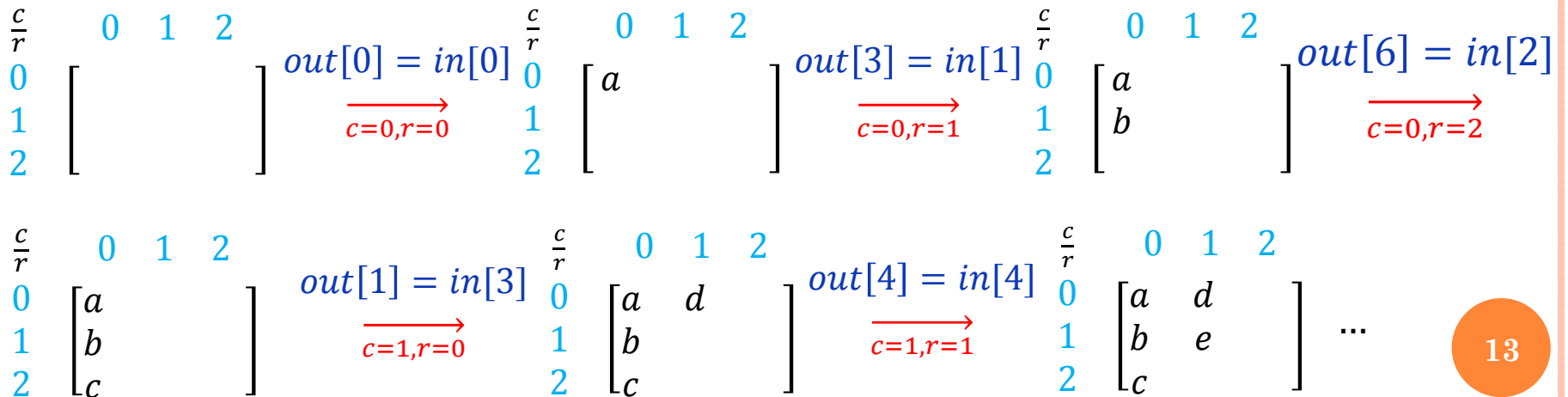
12

MATRIX TRANSPOSE - SERIAL

$$\begin{array}{c} \frac{c}{r} \\ 0 \\ 1 \\ 2 \end{array} \begin{array}{ccc} 0 & 1 & 2 \\ \left[\begin{array}{ccc} a & b & c \\ d & e & f \\ g & h & i \end{array} \right]
 \end{array} \rightarrow \begin{array}{cccccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ [a & b & c & d & e & f & g & h & i]
 \end{array}$$

```

for (int r = 0; r < N; r++)
    for (int c = 0; c < N; c++)
        Out[r][c] = in [c][r];
    
```



<Forward

MATRIX TRANSPOSE – PARALLEL

“SIMPLE”

```
for (int r = 0; r < N; r++)  
    for (int c = 0; c < N; c++)  
        out[r + c*N] = in[c + r*N]; // out(r,c) = in(c,r)
```

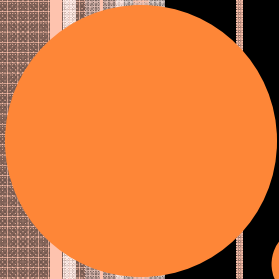
- # blocks=1, # threads= 1024

```
int r = threadIdx.x; //row by thread  
  
for (int c = 0; c < N; c++)  
    out[r + c*N] = in[c + r*N]; // out(r,c) = in(c,r)
```

- # blocks=1024, # threads= 1024

```
int r = blockIdx.x; //row  
int c = threadIdx.x; //column |  
  
out[r + c*N] = in[c + r*N]; // out(r,c) = in(c,r)
```



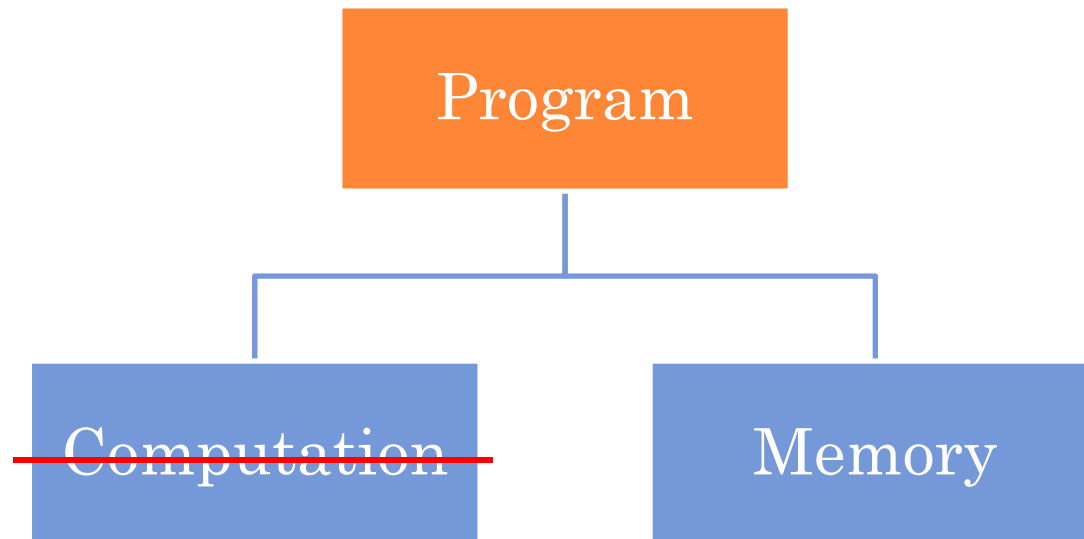


OPTIMIZE

Memory



OPTIMIZATION



Device Query (Start)

OPTIMIZATION (CONT.1)

- A: Memory clock: 900 Mhz = $900 * 10^6$ clock per sec.
- B: Memory bus bandwidth: 64-bit = 8 byte per clock

Quiz: What is the theoretical peak bandwidth?

$$= A * B = 7200 * 10^6 \text{ byte per sec.}$$

- Simple parallel algorithm:
 - N= 1024
 - A: Time= 6.70794 ms
 - B: Memory $1024 * 1024 * 4_{(\text{sizeof(int)})} * 2_{(\text{read and write})} = 2^{23}$
 - Bandwidth= B/A = $1250 * 10^6$ byte per sec
- Performance= $1250/7200 = 17\%$ <40 %
Bad

LITTLE'S LAW

$$\text{Useful bytes delivered} = \text{average latency} * \text{bandwidth}$$

OPTIMIZATION (CONT.3)

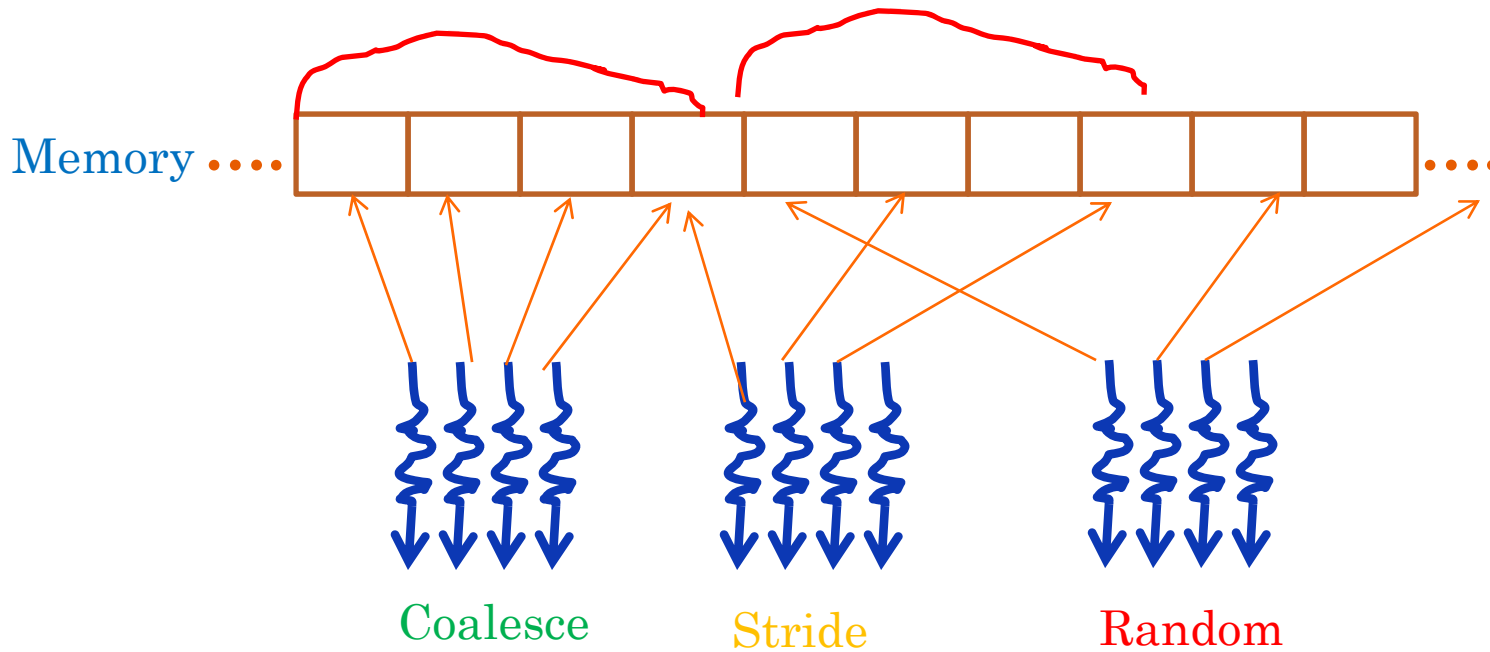
- nSight: Integrated debuggers and profilers, a complete development environment with many tools.
- Nvidia Visual Profiler (NVVP)

LITTLE'S LAW

$$\text{Useful bytes delivered} = \text{average latency} * \text{bandwidth}$$

- Make sure the moved data are useful

OPTIMIZATION (CONT.2)



```
int r = blockIdx.x; //row
int c = threadIdx.x; //column
out[r + c*N] = in[c + r*N]; // out(r,c) = in(c,r)
```

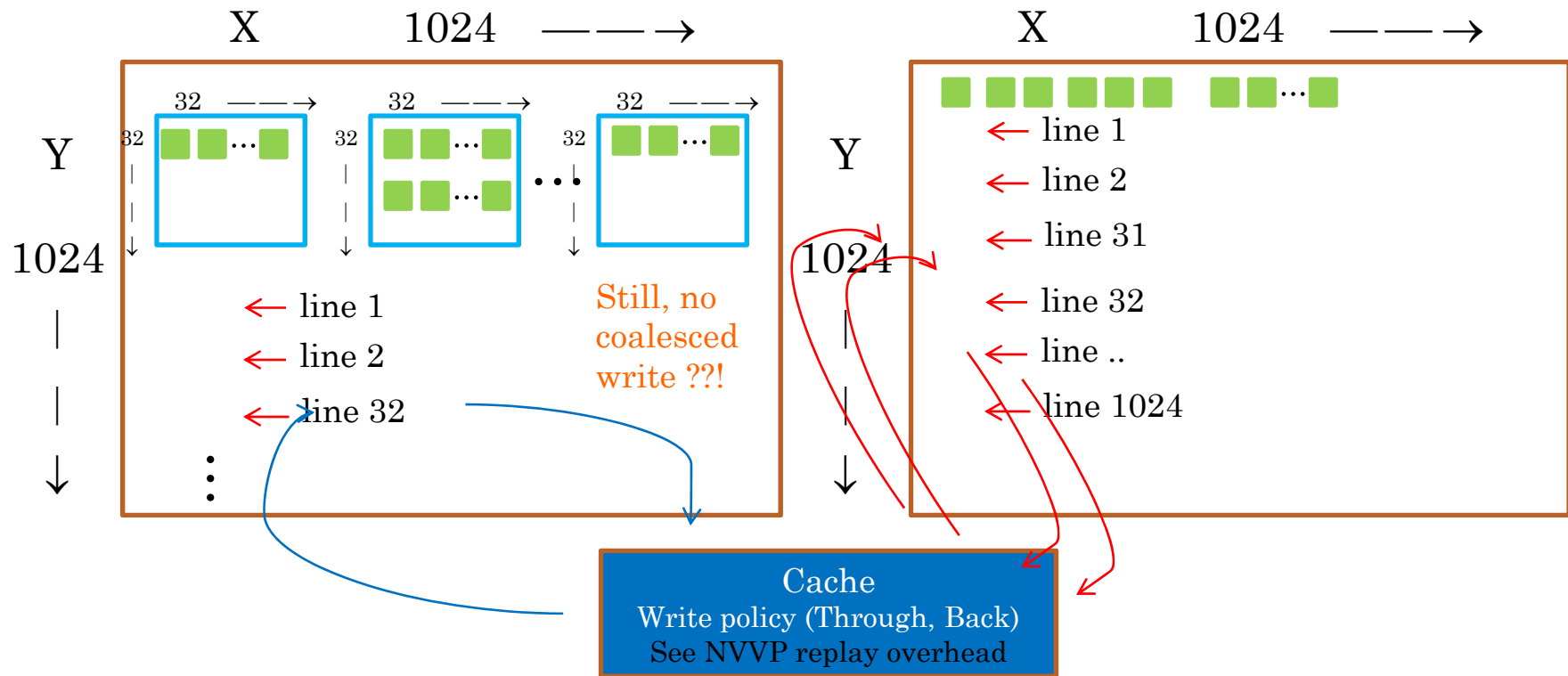
(N step)

(1 step)

[Back>](#)

MATRIX TRANSPOSE – PARALLEL “TILED” [COALESCED]

- # blocks=32*32, # threads= 32*32



```
int r = blockIdx.y * K + threadIdx.y;
int c = blockIdx.x * K + threadIdx.x;

out[r + c*N] = in[c + r*N]; // out(r,c) = in(c,r)
```



LITTLE'S LAW

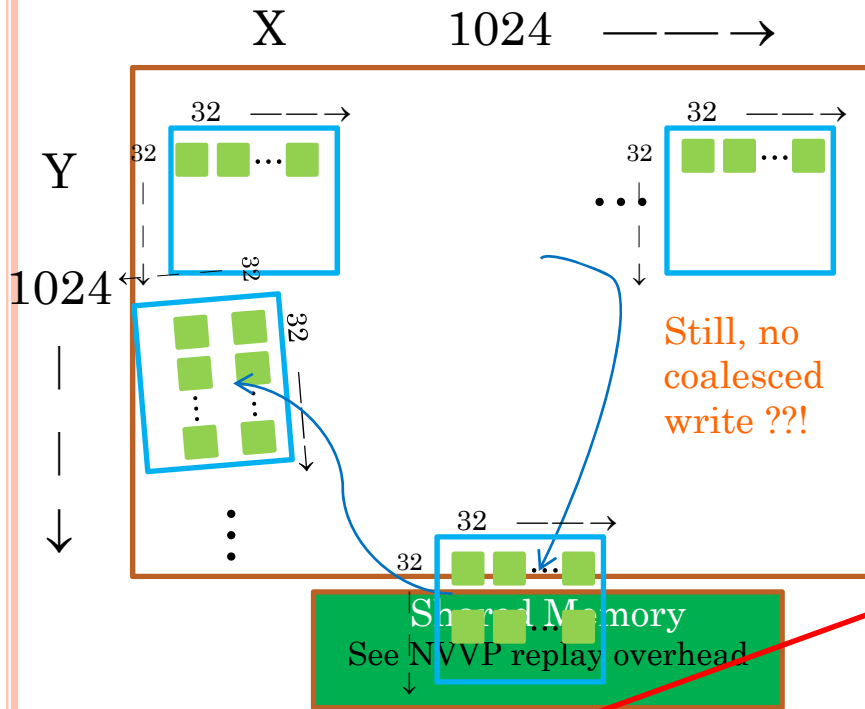
Useful bytes delivered = average  latency * bandwidth

- DRAM latency = 100's clock cycle
- Move to Shared memory

MATRIX TRANSPOSE – PARALLEL

“TILED SHARED MEMORY”

- # blocks=32*32, # threads= 32*32



Overhead

```
int x = threadIdx.x, y = threadIdx.y;
```

```
int in_r = blockIdx.y * K + y,
    out_r = blockIdx.x * K + y;
int in_c = blockIdx.x * K + x,
    out_c = blockIdx.y * K + x;
```

```
__shared__ float tile[K][K];
```

```
// coalesced read from global mem,
// write into shared mem:
tile[y][x] = in[in_c + in_r*N];
__syncthreads();
```

```
//TRANSPOSED read from shared mem,
//coalesced write to global mem:
out[out_c + out_r*N] = tile[x][y];
```

32 * 4 > warp limit
(See also NVVP replay overhead)

OCCUPANCY

- Each SM has a limited number of: Device Query (Start)
 - Blocks (currently 8)
 - Threads (currently 1536 to 2048)
 - Registers (currently 64K)
 - Shared memory (Currently 16k to 48k)
- All these factors are related:
 - Example 1:
 - If my GPU has 48 k shared memory
 - And my kernel takes 16k per block
 - Then my SM can only run 3 blocks at a time
 - Example 2:
 - If my SM can run 1536
 - And my kernel runs 1024 thread per block
 - Then my kernel can only run 1 block at a time
- If I run only 1024 thread then I waste 1536 – 1024 which considered bad Occupancy

OCCUPANCY (CONT.1)

- Check our tiled_shared kernel occupancy

Start

- How many blocks to run on an SM?

	My Kernel	My system Limit	Possible blocks
Registers/ Thread			
Shared memory			
Blocks/SM			

LITTLE'S LAW

Useful bytes delivered = average latency * bandwidth

- `_syncthreads()` cause overhead and decrease the latency because a thread may wait long time for other threads to finish.
- Decrease the number of threads per block



Make it 16*16
See NVVP Occupancy

OCCUPANCY (CONT.2)

- Shared memory bank conflict: shared memory organized in banks.

```
int x = threadIdx.x, y = threadIdx.y;
```

```
int in_r = blockIdx.y * K + y, out_r = blockIdx.x * K + y;
```

```
int in_c = blockIdx.x * K + x, out_c = blockIdx.y * K + x;
```

```
__shared__ float tile[K][K + 1];
```

```
// coalesced read from global mem, write into shared mem:
```

```
tile[y][x] = in[in_c + in_r*N];
```

```
__syncthreads();
```

```
// TRANSPOSED read from shared mem, coalesced write to global mem:
```

```
out[out_c + out_r*N] = tile[x][y];
```

The background features a black field with vertical stripes of varying widths and patterns. On the left, there are stripes with a fine grid pattern and a solid orange stripe. On the right, there is a solid orange stripe. Several orange circles of different sizes are scattered across the page, including a large one on the left and a medium one on the right.

29

OPTIMIZE Computation

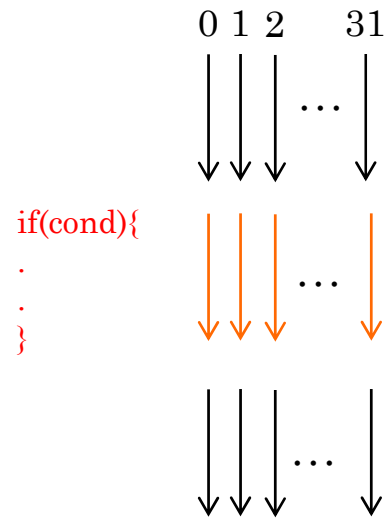
COMPUTATION OPTIMIZATION

- Minimize time spent at barriers
- Minimize thread divergence
- Math Optimization
- CPU-GPU interaction optimization
- Streams

MINIMIZE THREAD DIVERGENCE

- Warp: set of threads execute the same instruction at a time.
- SIMD: Single Instruction Multiple Data
 - CPU SSE / AVX vector registers
- SIMT: Single Instruction Multiple Thread

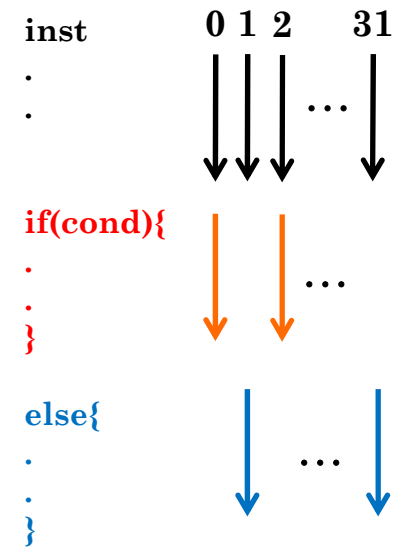
THREADS WARP



BRANCH DIVERGENCE

- What is the maximum branch divergence penalty for a CUDA thread block with 1024 threads?

32x slowdown



2-way branch divergence

QUIZ

What will be the slowdown for each of the following expressions in switch statement?

32 `__global__ void foo() { switch (threadIdx.x%32) case(0...31) }
kernel <<1024, 1>> ();`

32 `__global__ void foo() { switch (threadIdx.x%64) case(0...63) }
kernel <<1024, 1>> ();`

1 `__global__ void foo() { switch (threadIdx.y%32) case(0...31) } Next>
kernel <<64*16, 1>> ();`

2 `__global__ void foo() {switch (threadIdx.y%16) case(0...31) } Next>
kernel <<16*16, 1>> ();`

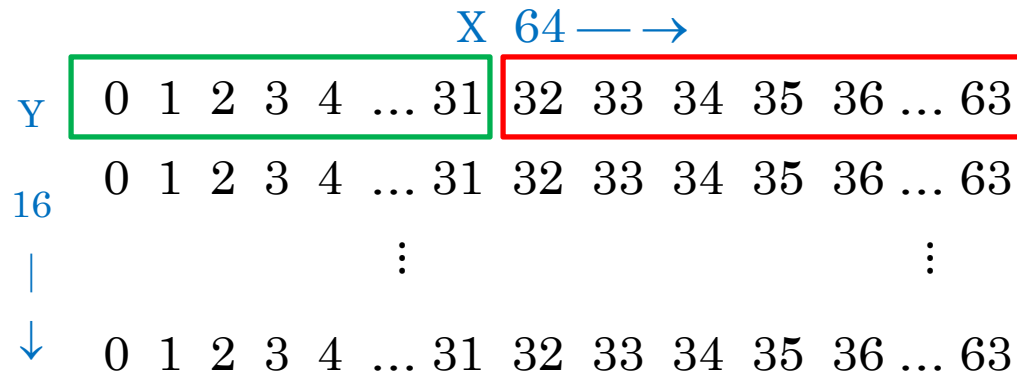
QUIZ (CONT.1)

- CUDA assigns thread ID as following:

x varies fast

y varies slower

z varies slowest



[<Back](#)

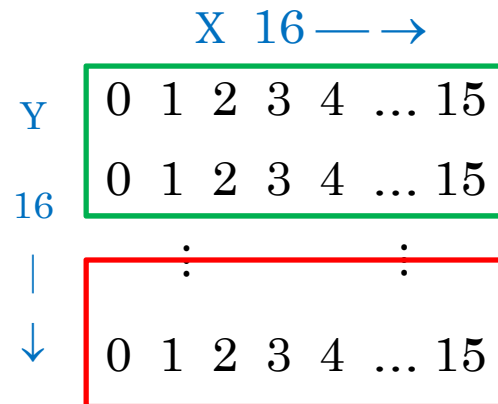
QUIZ (CONT.2)

- CUDA assigns thread ID as following:

x varies fast

y varies slower

z varies slowest



[<Back](#)

LOOP DIVERGENCE

```
__global__ void foo(){
    for(int i=0; i < threadIdx.x * 32; i++) //A
        bar();

    for(int i=0; i < threadIdx.x * 32; i++) //B
        bar();
}

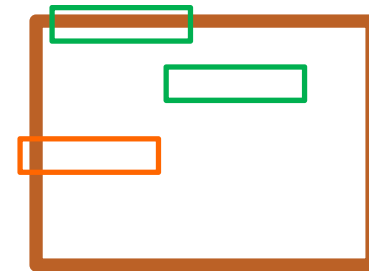
kernel <<1024, 1>> ();
```

- Which loop is faster?
 - 2nd
- How much fast?
 - 2 times

REAL WORLD DIVERGENCE EXAMPLE

- Operating on 1024 x 1024 image special handling of pixels on the boundary

```
__global__ perPixelKernel(float3 * image){  
    if(threadIdx.x==0 || threadIdx.x==1024 ||  
        threadIdx.y==0 || threadIdx.y==1024){  
        operation_on_boundries();  
    }  
    else{  
        operation_on_pixels();  
    }  
}
```



- What is the maximum branch divergence of any warp?
 - 2 way
- How many total threads are in diverged warps-warps with at least 2-way divergence?
 - $2 * 1022$

OPTIMIZE BRANCH DIVERGENCE

- Avoid branching code
 - Lots of if and switch statements
 - Try to restructure
- Beware large imbalance in thread workloads
 - Ex. When some threads perform the loop more than other threads

COMPUTATION OPTIMIZATION

- Minimize time spent at barriers
- Minimize thread divergence
- Math Optimization
- CPU-GPU interaction optimization
- Streams

MATH OPTIMIZATION

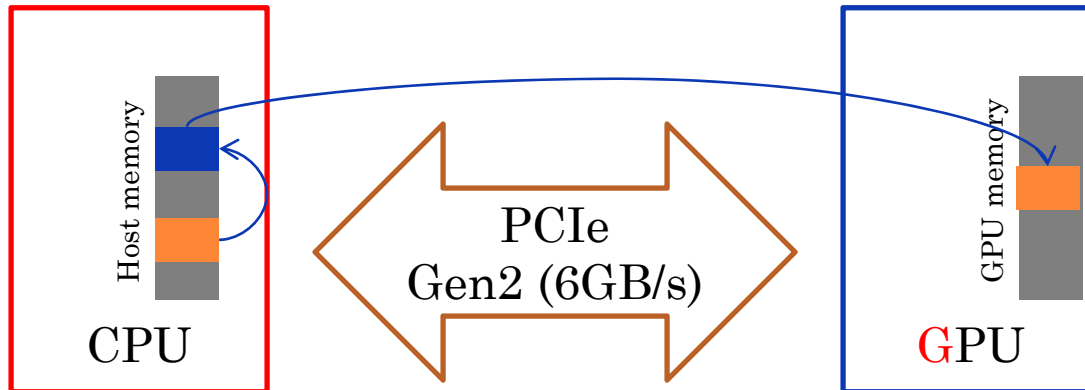
- Use double precision when only needed
- 64-bit math slower than 32-bit math
 - `float a=b+2.5 > float a=b+2.5f`
- Use intrinsics when possible
 - `__sin(), sqrt()` (2-3 bits **less precision** but **faster**)
`math.h`

COMPUTATION OPTIMIZATION

- Minimize time spent at barriers
- Minimize thread divergence
- Math Optimization
- CPU-GPU interaction optimization
- Streams

HOST-GPU INTERACTION

- Swap-out problem



HOST-GPU INTERACTION OPTIMIZATION

- Use pinned host memory (page-locked)
 - `cudaMallocHost`, `cudaHostRegister` instead of `Malloc`
- If you are using Pinned memory:
 - Use `cudaMemcpyAsync` instead of `cudaMemcpy`
 - This help to free the CPU
 - But you need to define stream

COMPUTATION OPTIMIZATION

- Minimize time spent at barriers
- Minimize thread divergence
- Math Optimization
- CPU-GPU interaction optimization
- Streams

STREAMS

- Stream: sequence of operations that execute in order.
 - (Memory transfers, kernels)

- Ex.

```
cudaMemcpyAsync(..);      cudaMemcpyAsync(.., s1);  
A<<..>>;                 A<<>>(.., s2);  
cudaMemcpyAsync(..);      cudaMemcpyAsync(.., s3);  
B<<..>>; ✗                B<<>>(.., s4); ✓
```

Default steam (s0)

```
cuda stream_t s1;  
cuda streamCreate(s1);  
..  
cuda streamDestroy(s1);
```

QUIZ

- If all operations take exactly 1 second, how long to complete the following:

```
cudaStream_t s1, s2;
cudaStreamCreate(&s1); cudaStreamCreate(&s2);

3 {
  cudaMemcpy(&d_arr, &h_arr, numbytes, cudaH2D);
  A<<<1, 128>>>(d_arr);
  cudaMemcpy(&h_arr, &d_arr, numbytes, cudaD2H);
}

3 {
  cudaMemcpyAsync(&d_arr, &h_arr, numbytes, cudaH2D, s1);
  A<<<1, 128, s1>>>(d_arr);
  cudaMemcpyAsync(&h_arr, &d_arr, numbytes, cudaD2H, s1);
}

3 {
  cudaMemcpyAsync(&d_arr1, &h_arr1, numbytes, cudaH2D, s1);
  A<<<1, 128, s1>>>(d_arr1);
  cudaMemcpyAsync(&h_arr1, &d_arr1, numbytes, cudaD2H, s1);
  cudaMemcpyAsync(&d_arr2, &h_arr2, numbytes, cudaH2D, s2);
  B<<<1, 192, s2>>>(d_arr2);
  cudaMemcpyAsync(&h_arr2, &d_arr2, numbytes, cudaD2H, s2);
}

3 {
  cudaMemcpyAsync(&d_arr1, &h_arr1, numbytes, cudaH2D, s1);
  cudaMemcpyAsync(&d_arr2, &h_arr2, numbytes, cudaH2D, s2);
  A<<<1, 128, s1>>>(d_arr1);
  B<<<1, 192, s2>>>(d_arr2);
  cudaMemcpyAsync(&h_arr1, &d_arr1, numbytes, cudaD2H, s1);
  cudaMemcpyAsync(&h_arr2, &d_arr2, numbytes, cudaD2H, s2);
}
```

QUIZ [CONT.]

- If all operations take exactly 1 second, how long to complete the following:

```
cudaStream_t s1, s2;  
cudaStreamCreate(&s1); cudaStreamCreate(&s2);
```

1

```
cudaMemcpyAsync(&d_arr1, &h_arr1, numbytes, cudaH2D, s1);  
A<<<1, 128, s2>>>(d_arr2);
```

1

```
cudaMemcpyAsync(&d_arr1, &h_arr1, numbytes, cudaH2D, s1);  
A<<<1, 128, s2>>>(d_arr1);
```

With
a wrong values

STREAMS (CONT.)

- Old GPUs big amount of data to chunks
 - Too much overhead to transfer data back and forward

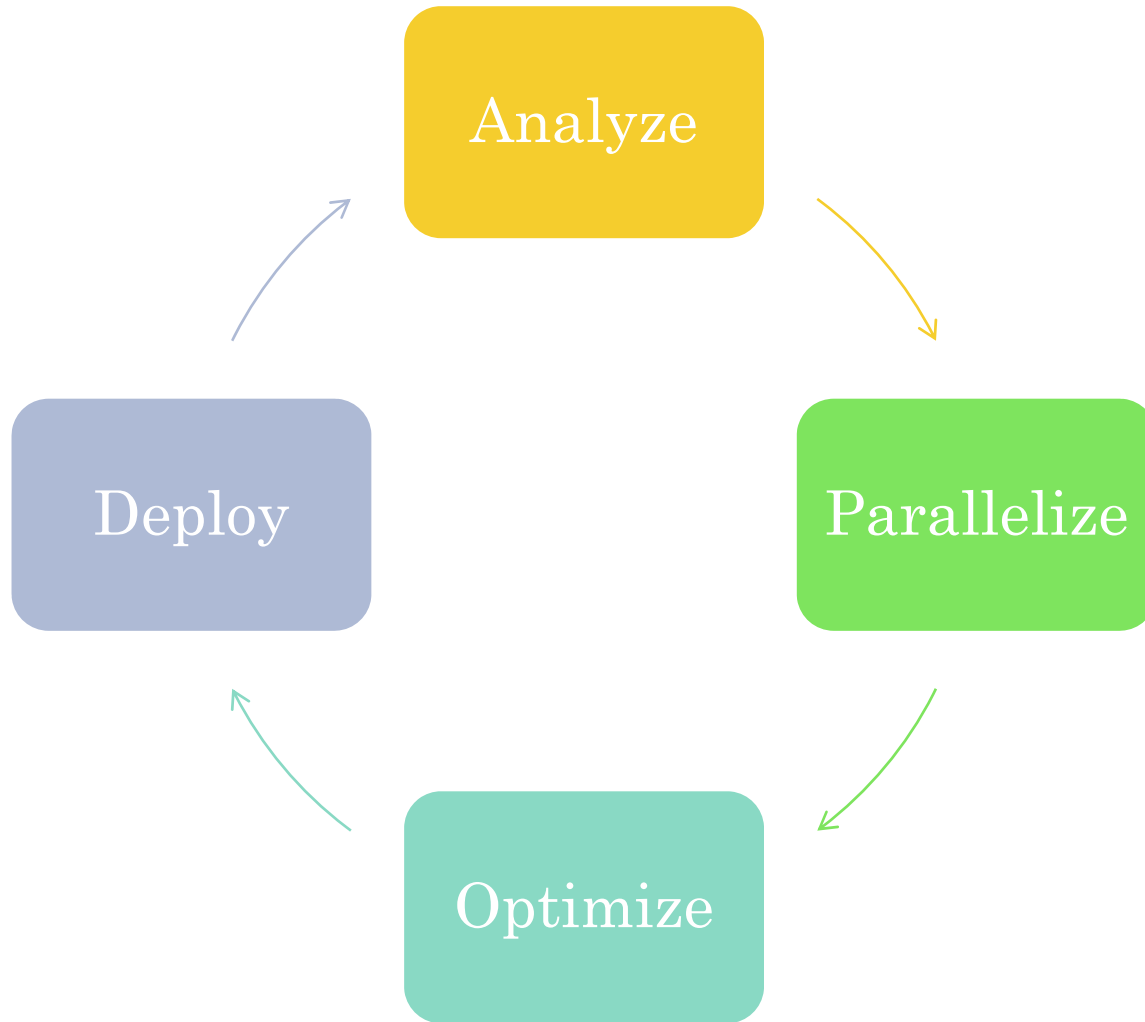


- Streams good with big amount of data
 - Run chunks and kernels in asynchronous manner
 - Overlap memory and compute
 - Fill GPU with smaller kernels

COMPUTATION OPTIMIZATION

- Minimize time spent at barriers
- Minimize thread divergence
- Math Optimization
- CPU-GPU interaction optimization
- Streams

APOD



PROBLEM SET

- Optimize the Histogram program in lecture 3
 - Use local bins instead of atomicAdd

