



PROFILING & OPTIMIZING GPU KERNELS

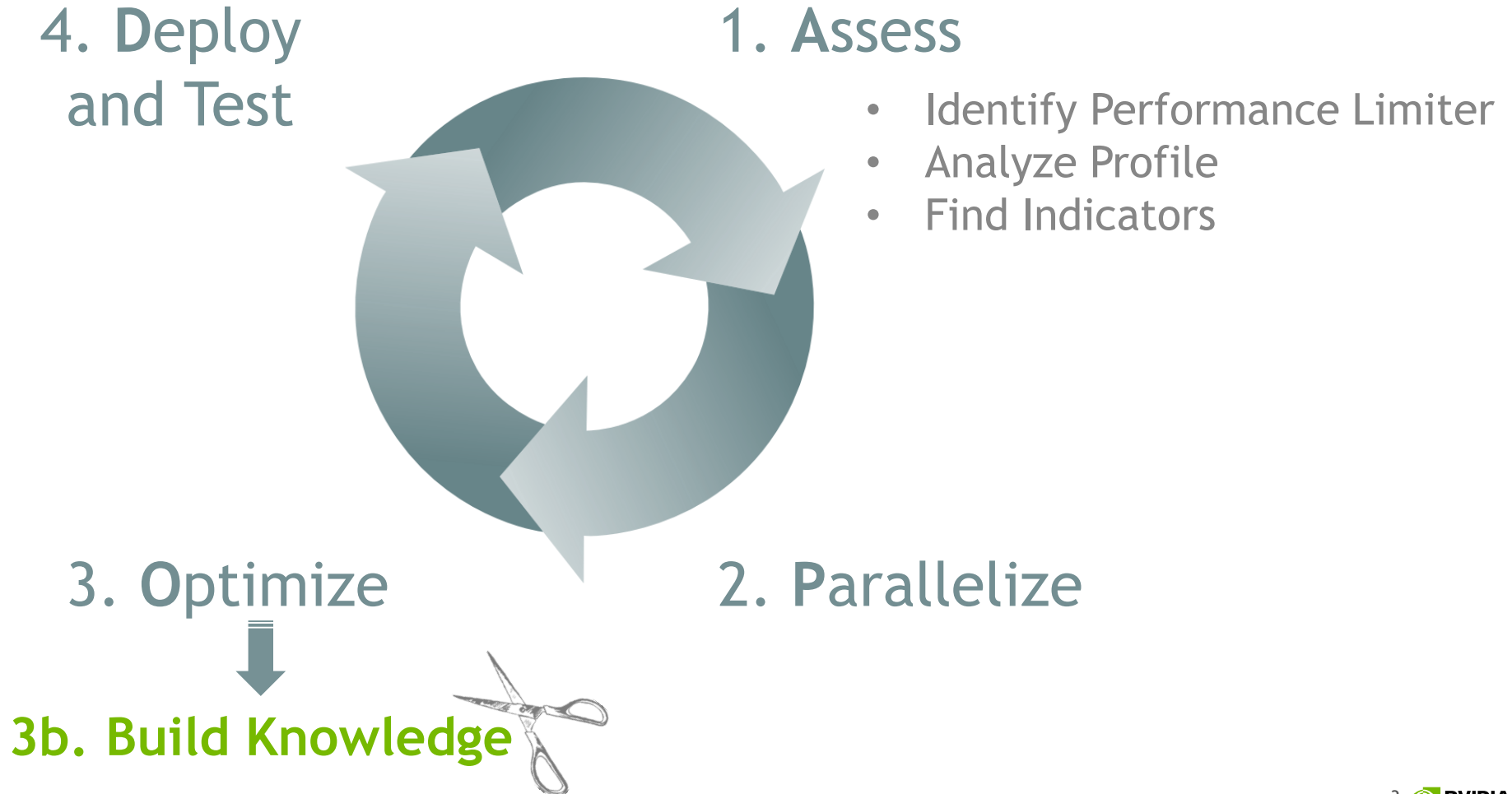
Vishal Mehta

BEFORE YOU START

The five steps to enlightenment

1. Know your application
 - What does it compute? How is it parallelized? What final performance is expected?
2. Know your hardware
 - What are the target machines, how many nodes? Machine-specific optimizations okay?
3. Know your tools
 - Strengths and weaknesses of each tool? Learn how to use them (and learn one well!)
4. Know your process
 - Performance optimization is a constant learning process
5. Make it so!

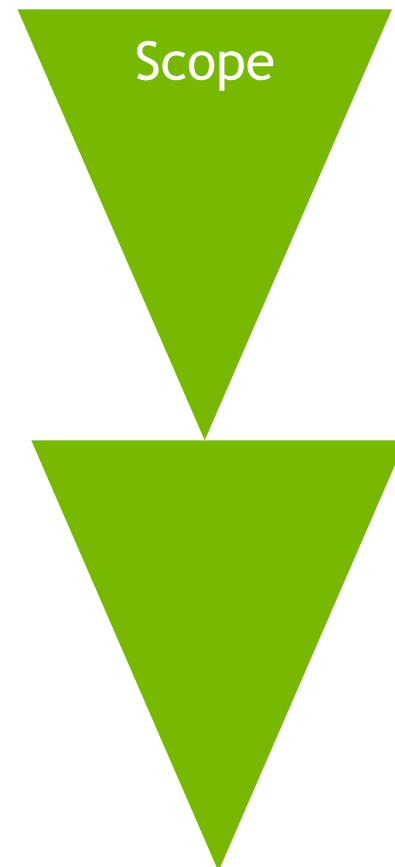
THE APOD CYCLE



GUIDING OPTIMIZATION EFFORT

“Drilling Down into the Metrics”

- Challenge: How to know where to start?
- Top-down Approach:
 - Find Hotspot Kernel
 - Identify Performance Limiter of the Hotspot
 - Find performance bottleneck indicators related to the limiter
 - Identify associated regions in the source code
 - Come up with strategy to fix and change the code
 - Start again

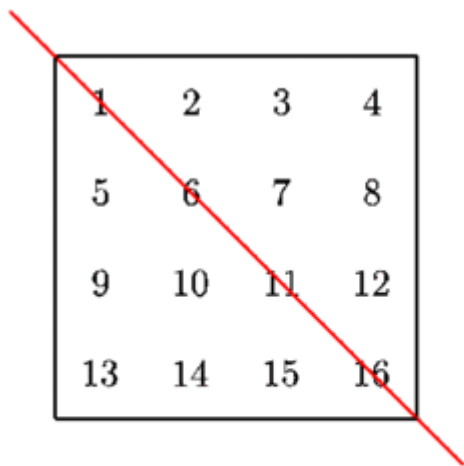


**KNOW YOUR APPLICATION:
MATRIX TRANSPOSE**

MATRIX TRANSPOSE

Parallel transpose kernel

- Motivation: Used in lot of Linear Algebra, Relatively easy for hands-on, but still a lot that could be done to optimize



1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

A

```
// Simplest transpose; doesn't use shared memory.  
// Global memory reads are coalesced but writes are not.  
__global__ void transposeNaive(float *odata, const float *idata)  
{  
    int x = blockIdx.x * TILE_DIM + threadIdx.x;  
    int y = blockIdx.y * TILE_DIM + threadIdx.y;  
    int width = gridDim.x * TILE_DIM;  
  
    for (int j = 0; j < TILE_DIM; j += BLOCK_ROWS)  
        odata[x*width + (y+j)] = idata[(y+j)*width + x];  
}
```

KNOW YOUR HARDWARE: PASCAL ARCHITECTURE

GPU COMPARISON

	V100	P100 (SXM2)	M40	K40
Double/Single/Half TFlop/s	7.8/15.7/125 (TensorCore)	5.3/10.6/21.2	0.2/7.0/NA	1.4/4.3/NA
Memory Bandwidth (GB/s)	900	732	288	288
Memory Size	16/32 GB	16GB	12GB, 24GB	12GB
L2 Cache Size	6144 KB	4096 KB	3072 KB	1536 KB
Base/Boost Clock (Mhz)	1312 / 1530	1328/1480	948/1114	745/875
TDP (Watts)	300	300	250	235

GP100 SM

GP100

CUDA Cores 64

Register File 256 KB

Shared Memory 64 KB

Active Threads 2048

Active Blocks 32



KNOW YOUR TOOLS: PROFILERS

PROFILING TOOLS

Many Options!

From NVIDIA

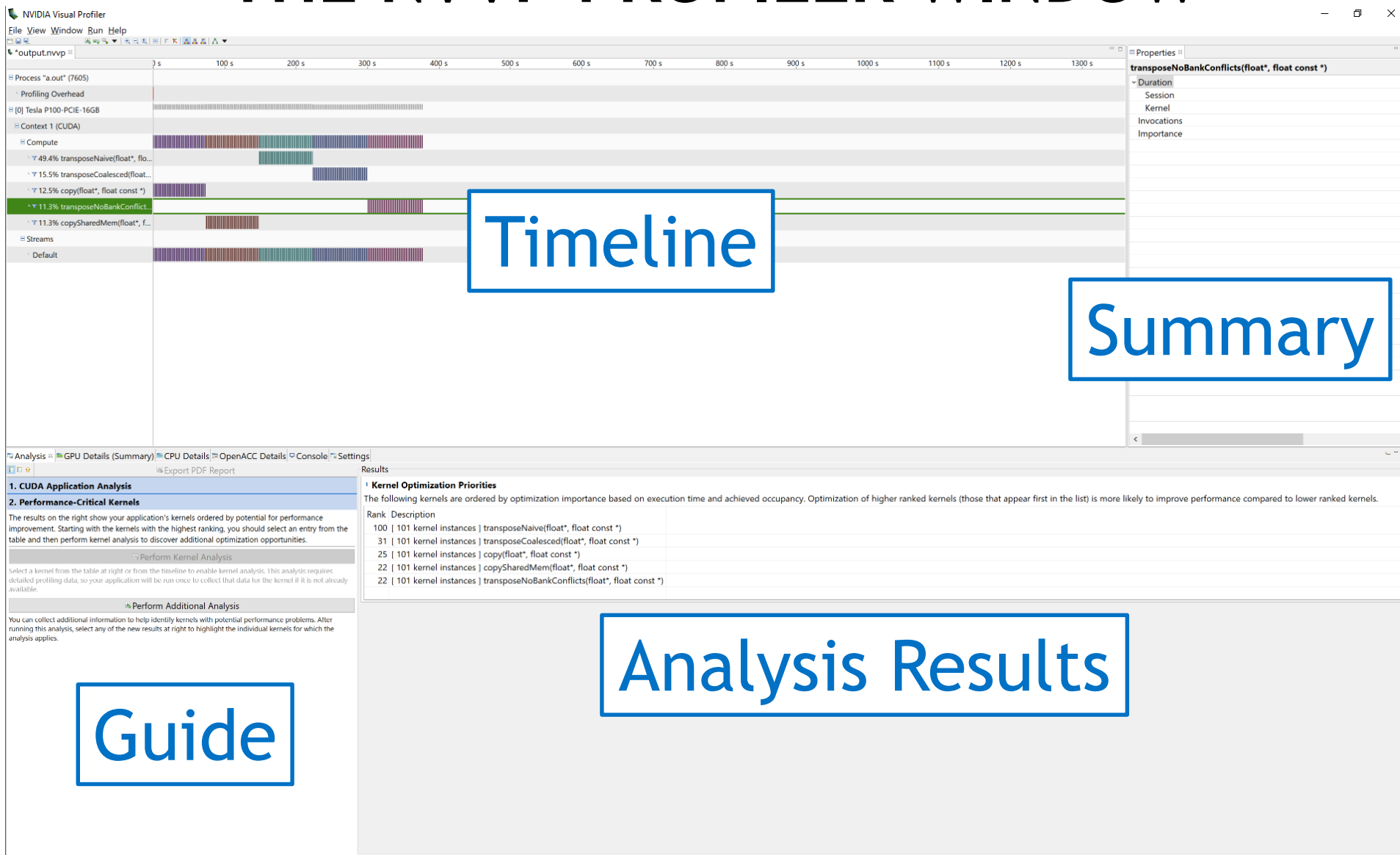
- nvprof
- NVIDIA Visual Profiler
 - Standalone (nvvp)
 - Integrated into Nsight Eclipse Edition (nsight)
- Nsight Visual Studio Edition

Third Party

- TAU Performance System
- VampirTrace
- PAPI CUDA component
- HPC Toolkit
- (Tools using CUPTI)

In this session we will be showing nvvp screenshots

THE NVVP PROFILER WINDOW



**MAKE IT HAPPEN:
ITERATION 1
GLOBAL MEMORY COALESCING**

IDENTIFY HOTSPOT

Hotspot



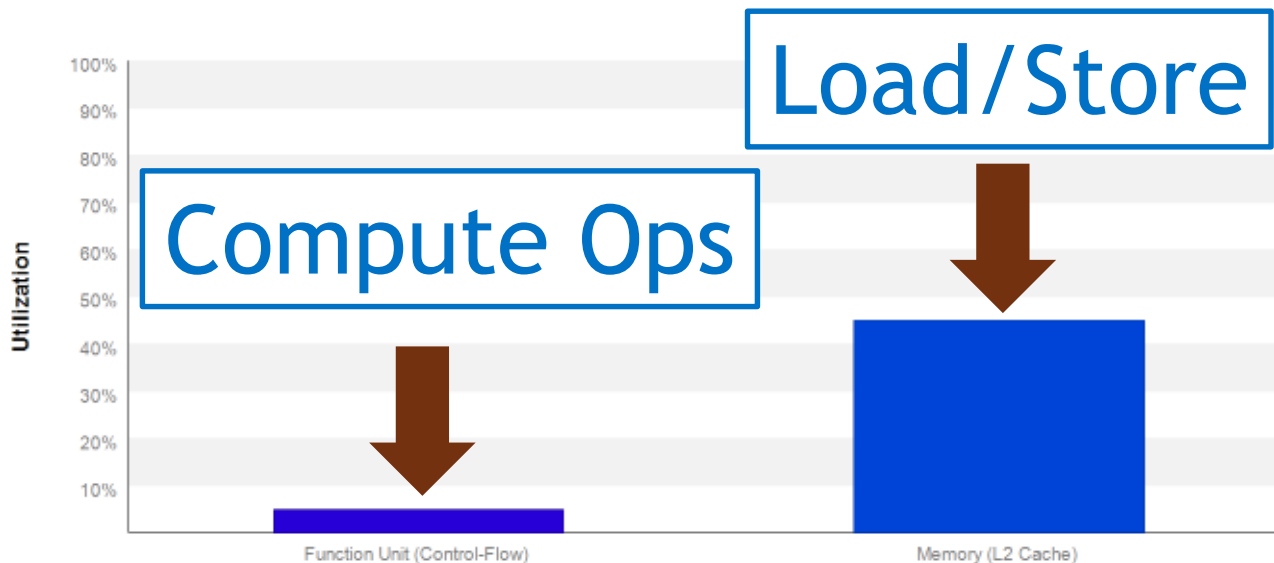
Results		
Kernel Optimization Priorities		
The following kernels are ordered by optimization importance based on execution time and achieved occupancy. Optimization of higher ranked kernels (those that appear first in the list) is more likely to improve performance compared to lower ranked kernels.		
Rank	Description	
100	[101 kernel instances] transposeNaive(float*, float const *)	
31	[101 kernel instances] transposeCoalesced(float*, float const *)	
25	[101 kernel instances] copy(float*, float const *)	
22	[101 kernel instances] copySharedMem(float*, float const *)	
22	[101 kernel instances] transposeNoBankConflicts(float*, float const *)	

IDENTIFY PERFORMANCE LIMITER

Results

Kernel Performance Is Bound By Instruction And Memory Latency

This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of "Tesla P100-PCIE-16GB". These utilization levels indicate that the performance of the kernel is most likely limited by the latency of arithmetic or memory operations. Achieved compute throughput and/or memory bandwidth below 60% of peak typically indicates latency issues.

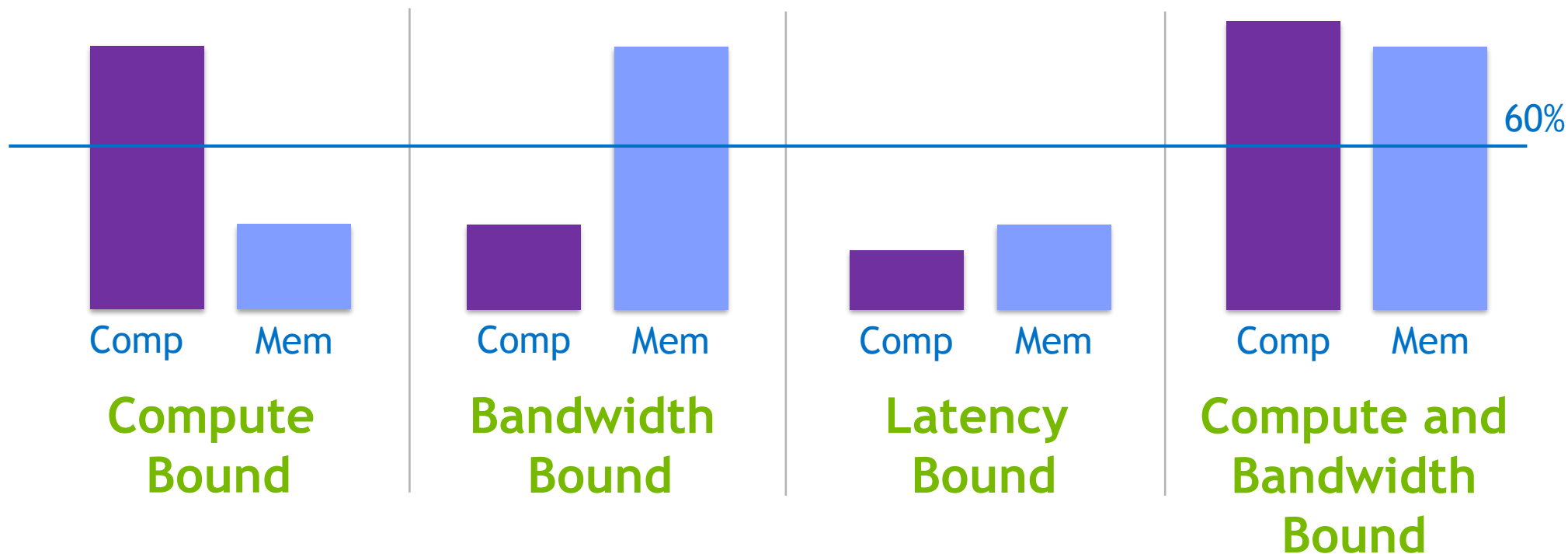


Memory Utilization Issues?

PERFORMANCE LIMITER CATEGORIES

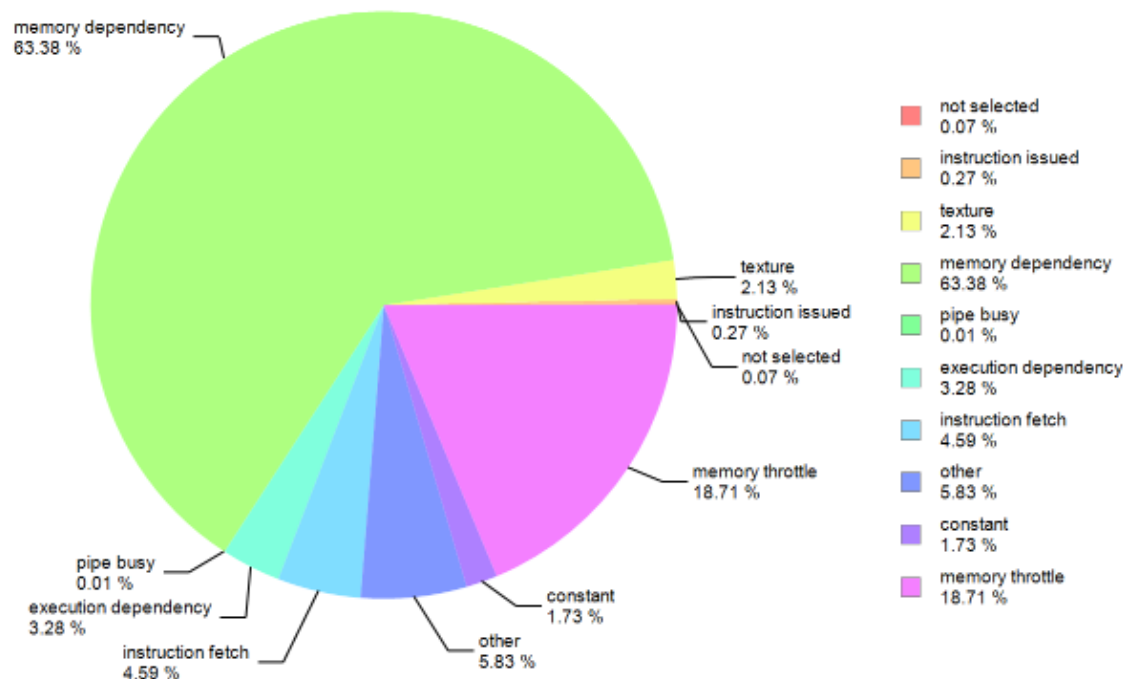
Memory Utilization vs Compute Utilization

Four possible combinations:



DRILLING DOWN: LATENCY ANALYSIS

Grid Size	[65536,1,1]
Block Size	[8,4,1]
Occupancy	
Achieved	⚠ 49.7%
Theoretical	50%
Limiter	Block Size



⚠ GPU Utilization May Be Limited By Block Size

Theoretical occupancy is less than 100% but is large enough that increasing occupancy may not improve performance. You can attempt the following optimization to increase the number of warps on each SM but it may not lead to increased performance.

OCCUPANCY

GPU Utilization

Each SM has limited resources:

- max. 64K Registers (32 bit) distributed between threads
- max. 48KB of shared memory per block (96KB per SMM)
- max. 32 Active Blocks per SMM
- Full occupancy: 2048 threads per SM (64 warps)

When a resource is used up, occupancy is reduced

() Values vary with Compute Capability*

LOOKING FOR MORE INDICATORS

The screenshot displays the NVIDIA Nsight Systems interface. The top-left pane shows the source code for `transposeNaive` in `transpose_result.cu`. Line 105 is highlighted in red. The top-right pane shows the corresponding disassembly, with instructions `STGE [R2], R11`, `LDGE R0, [R4]`, `STGE [R2+0x20], R0`, `LDGE R7, [R6]`, `STGE [R2+0x40], R7`, and `LDGE R8, [R8]` highlighted in blue. A large blue arrow points from the disassembly pane to the source code pane, labeled "Source Code Association". The bottom pane shows the "Results" section, specifically the "Global Memory Alignment and Access Pattern" analysis. It provides a summary of memory bandwidth usage and a table of results for line 105.

Line / File	transpose_result.cu - \users\vmeh\summerschool
105	Global Store L2 Transactions/Access = 32, Ideal Transactions/Access = 4 [262144 L2 transactions for 8192 total executions]
105	Global Store L2 Transactions/Access = 32, Ideal Transactions/Access = 4 [262144 L2 transactions for 8192 total executions]
105	Global Store L2 Transactions/Access = 32, Ideal Transactions/Access = 4 [262144 L2 transactions for 8192 total executions]
105	Global Store L2 Transactions/Access = 32, Ideal Transactions/Access = 4 [262144 L2 transactions for 8192 total executions]

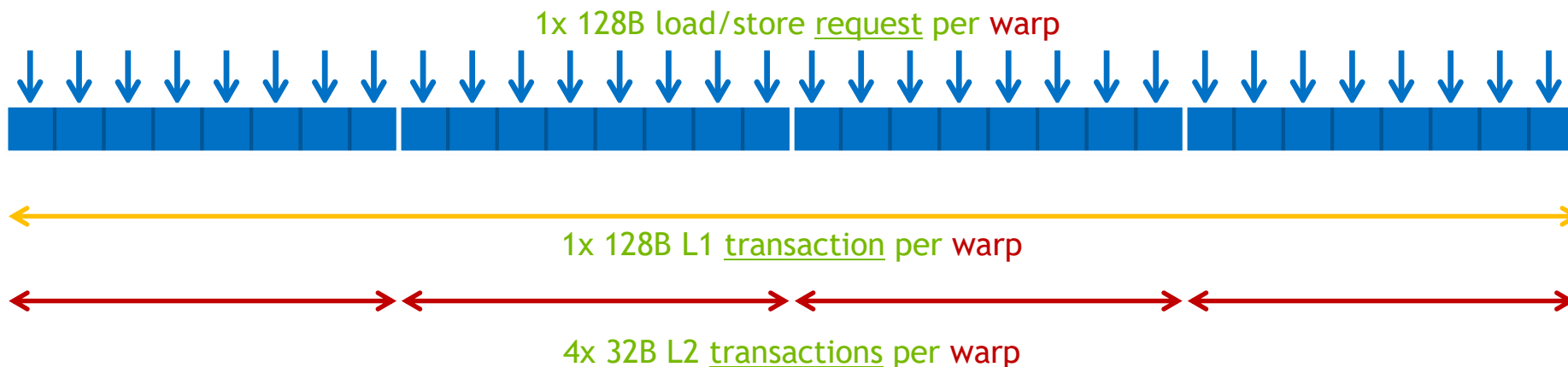
For line numbers use:
`nvcc -lineinfo`

32 Global STORE
Transactions

MEMORY TRANSACTIONS: BEST CASE

A warp issues 32x4B aligned and consecutive load/store request

Threads read different elements of the same 128B segment



1x L1 transaction: 128B needed / 128B transferred

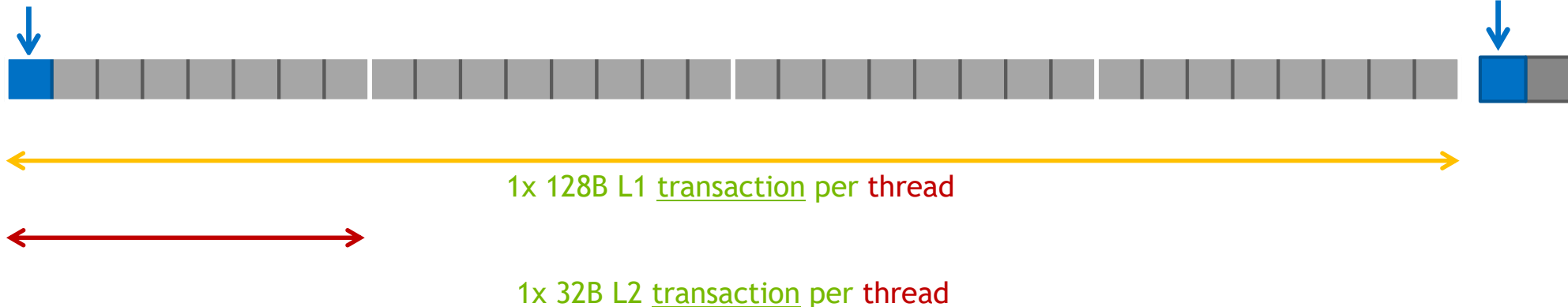
4x L2 transactions: 128B needed / 128B transferred

MEMORY TRANSACTIONS: WORST CASE

Threads in a warp read/write 4B words, 128B between words

Each thread reads the first 4B of a 128B segment

Stride: 32x4B



32x L1 transactions: 128B needed / **32x** 128B transferred

32x L2 transactions: 128B needed / **32x** 32B transferred

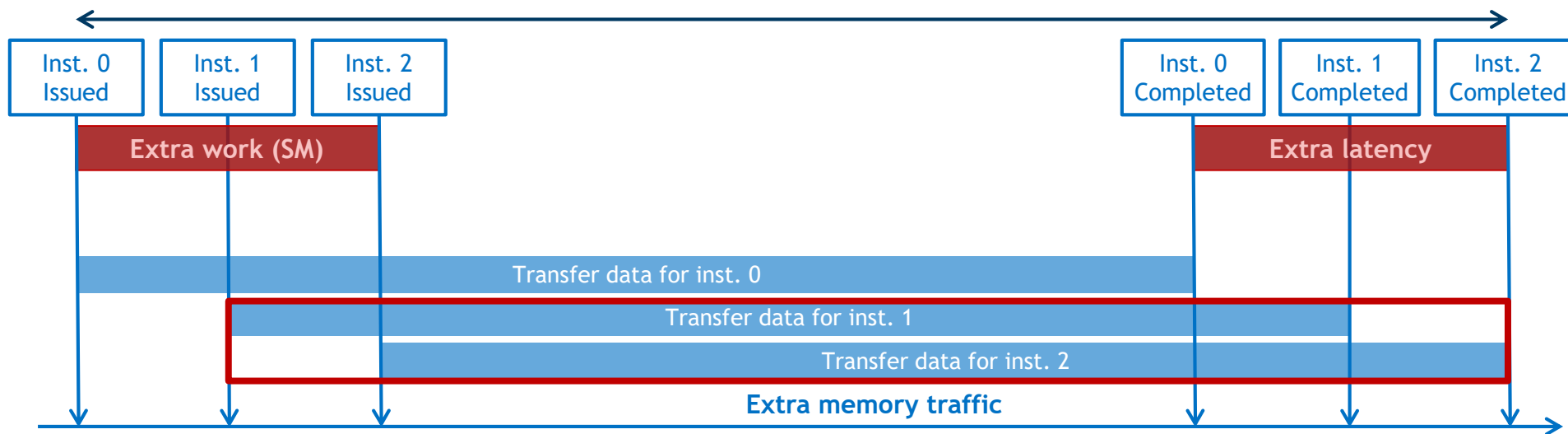
TRANSACTIONS AND REPLAYS

With replays, requests take more time and use more resources

More instructions issued

More memory traffic

Increased execution time



PERF-OPT QUICK REFERENCE CARD

Category:	Latency Bound - Occupancy
Problem:	Latency is exposed due to low occupancy
Goal:	<u>Hide</u> latency behind more parallel work
Indicators:	Occupancy low (< 60%) Execution Dependency High
Strategy:	Increase occupancy by: <ul style="list-style-type: none">• Varying block size• Varying shared memory usage• Varying register count (use <code>__launch_bounds</code>)



PERF-OPT QUICK REFERENCE CARD

Category:	Latency Bound - Coalescing
Problem:	Memory is accessed inefficiently => high latency
Goal:	Reduce #transactions/request to reduce latency
Indicators:	Low global load/store efficiency, High #transactions/#request compared to ideal
Strategy:	Improve memory coalescing by: <ul style="list-style-type: none">• Cooperative loading inside a block• Change block layout• Aligning data• Changing data layout to improve locality



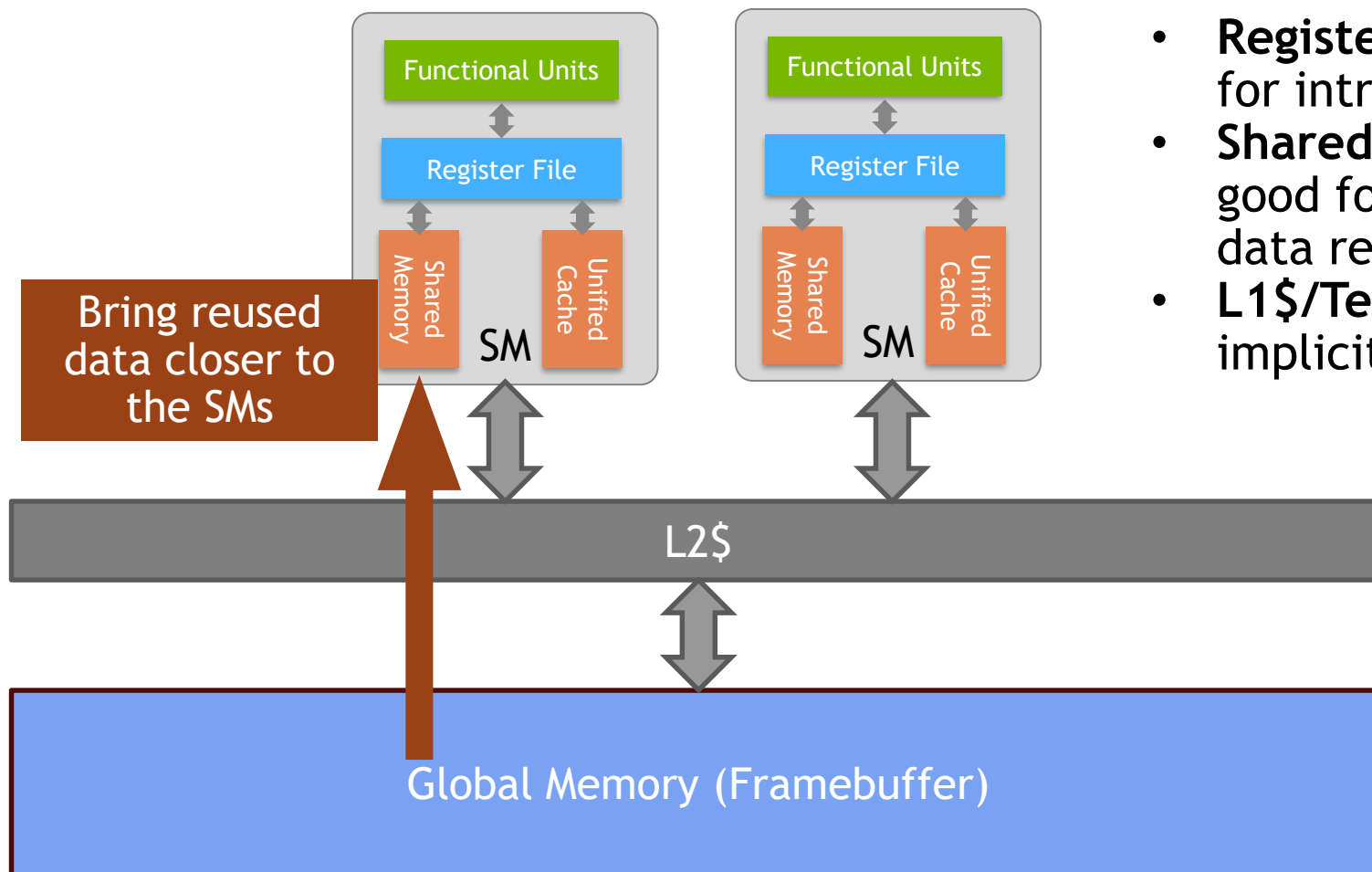
PERF-OPT QUICK REFERENCE CARD

Category:	Bandwidth Bound - Coalescing
Problem:	Too much unused data clogging memory system
Goal:	Reduce traffic, move more <u>useful</u> data per request
Indicators:	Low global load/store efficiency, High #transactions/#request compared to ideal
Strategy:	Improve memory coalescing by: <ul style="list-style-type: none">• Cooperative loading inside a block• Change block layout• Aligning data• Changing data layout to improve locality



GPU MEMORY HIERARCHY

P100 (SMX2)

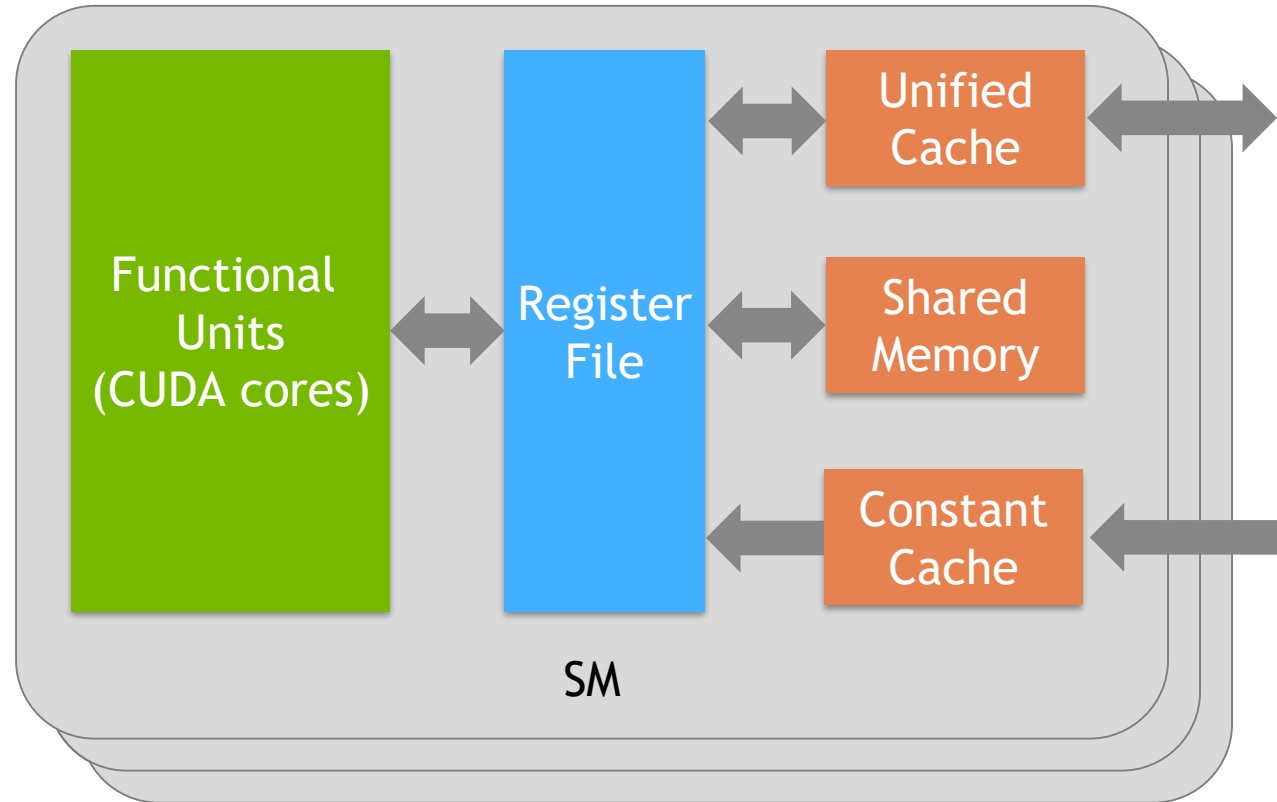


- **Registers (256 KB/SM):** good for intra-thread data reuse
- **Shared memory (64 KB/SM):** good for explicit intra-block data reuse
- **L1\$/Tex\$, L2\$ (4096 KB):** implicit data reuse

GPU SM ARCHITECTURE

Pascal SM

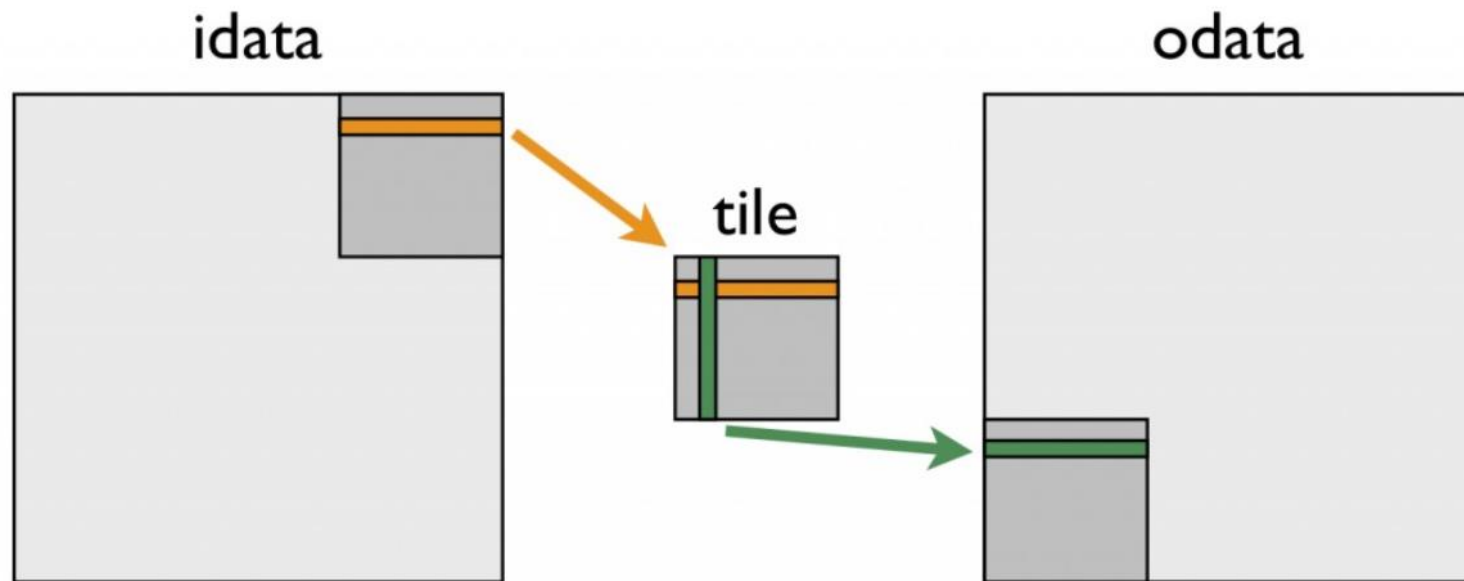
GP100	
CUDA Cores	64
Register File	256 KB
Shared Memory	64 KB



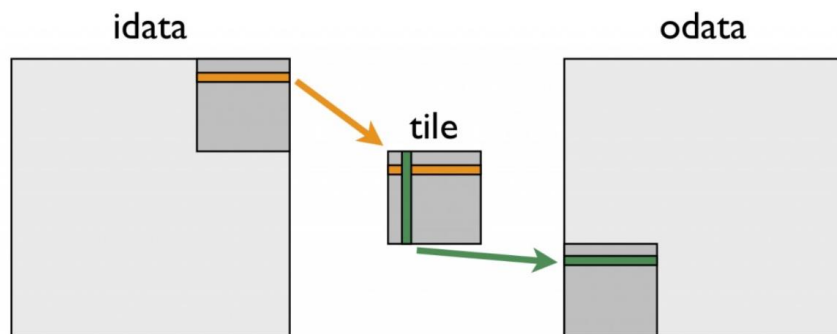
56 SMs on Tesla P100

USING SHARED MEMORY IN MATRIX TRANSPOSE

TILED MATRIX TRANSPOSE



TILED MATRIX TRANSPOSE



Device : Tesla P100-PCIE-16GB
Matrix size: 1024 1024, Block size: 32 8, Tile size: 32 32
dimGrid: 32 32 1. dimBlock: 32 8 1

Routine	Bandwidth (GB/s)
copy	442.94
shared memory copy	472.60
naive transpose	111.13
coalesced transpose	363.88

```
// coalesced transpose
// Uses shared memory to achieve coalescing in both reads and writes
// Tile width == #banks causes shared memory bank conflicts.
__global__ void transposeCoalesced(float *odata, const float *idata)
{
    __shared__ float tile[TILE_DIM][TILE_DIM];

    int x = blockIdx.x * TILE_DIM + threadIdx.x;
    int y = blockIdx.y * TILE_DIM + threadIdx.y;
    int width = gridDim.x * TILE_DIM;

    for (int j = 0; j < TILE_DIM; j += BLOCK_ROWS)
        tile[threadIdx.y+j][threadIdx.x] = idata[(y+j)*width + x];

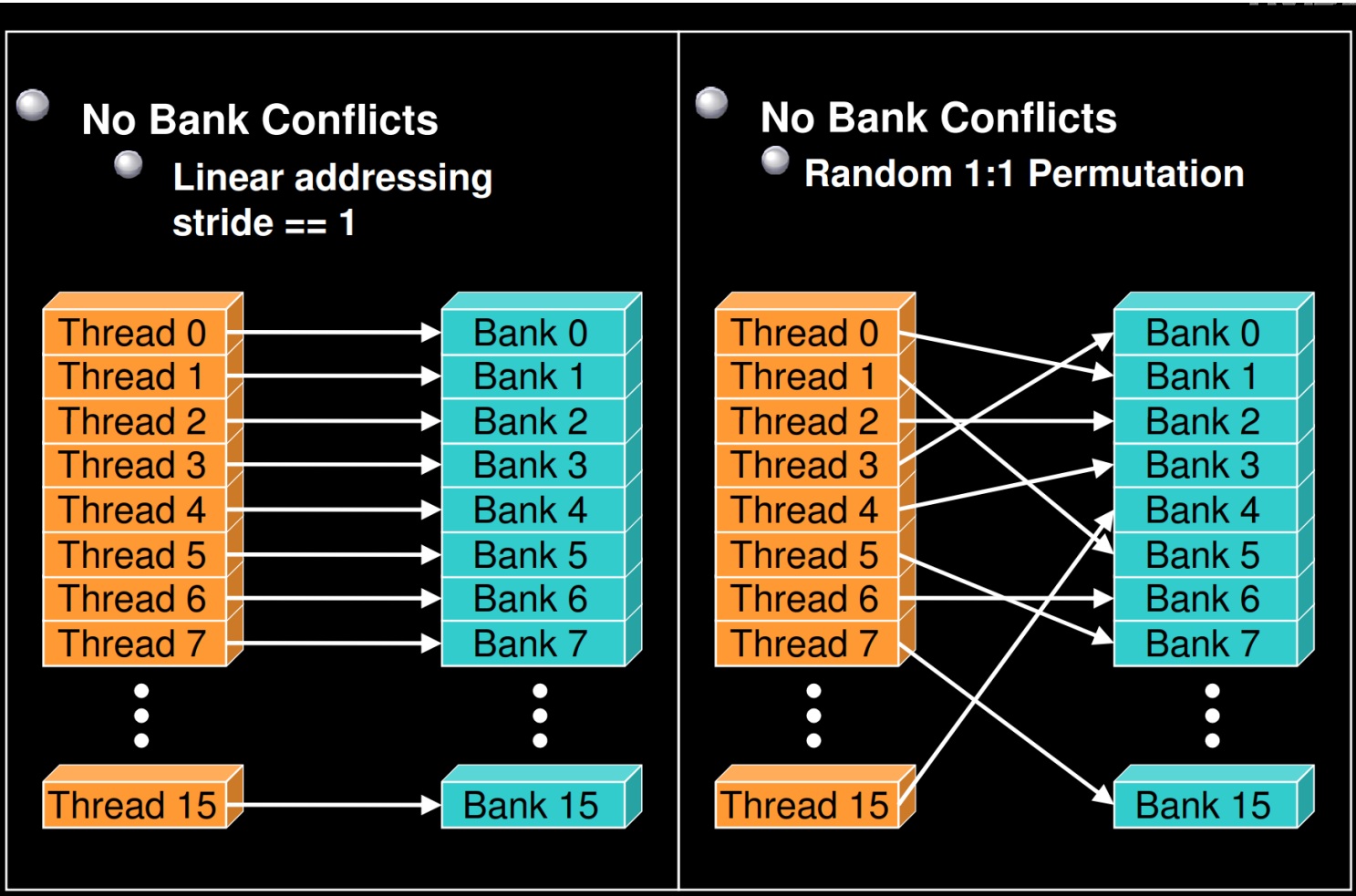
    __syncthreads();

    x = blockIdx.y * TILE_DIM + threadIdx.x; // transpose block offset
    y = blockIdx.x * TILE_DIM + threadIdx.y;

    for (int j = 0; j < TILE_DIM; j += BLOCK_ROWS)
        odata[(y+j)*width + x] = tile[threadIdx.x][threadIdx.y + j];
}
```

ITERATION 2: SHARED MEMORY OPTIMIZATION

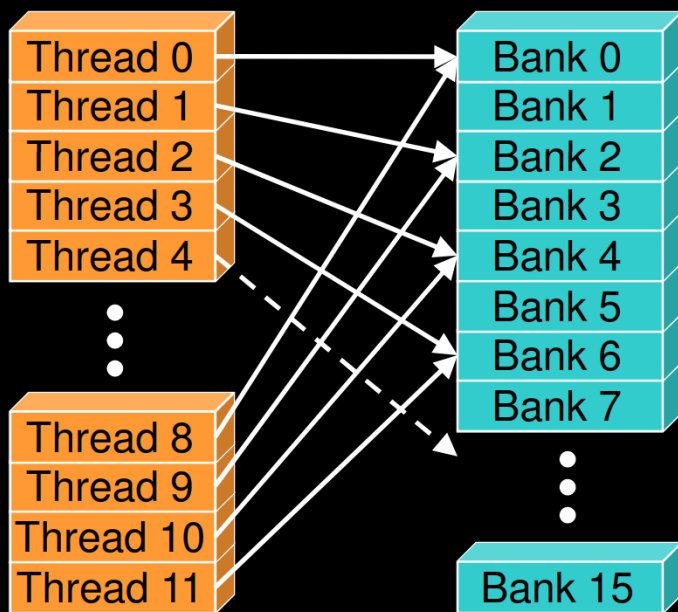
SHARED MEMORY EXAMPLES



SHARED MEMORY EXAMPLES

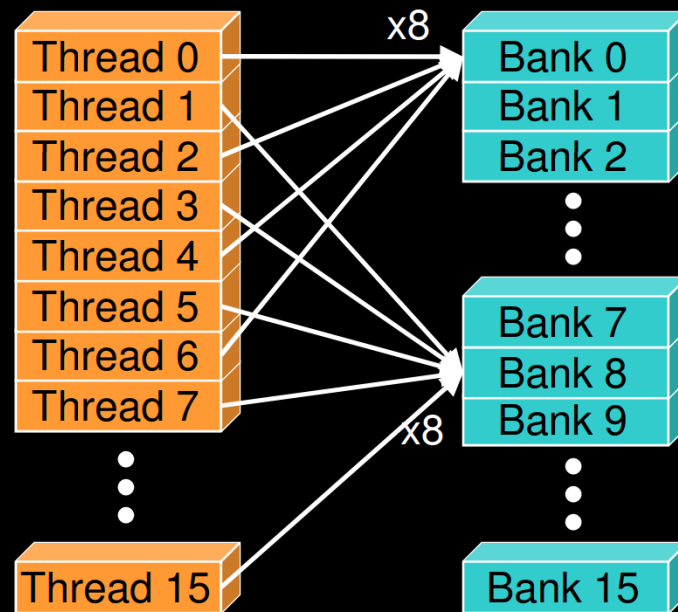
- 2-way Bank Conflicts

- Linear addressing stride == 2



- 8-way Bank Conflicts

- Linear addressing stride == 8



PERF-OPT QUICK REFERENCE CARD

Category:	Device Mem Bandwidth Bound - Shared Memory
Problem:	Too much data movement
Goal:	<u>Reduce</u> amount of data traffic to/from global mem
Indicators:	Higher than expected memory traffic to/from global memory Low arithmetic intensity of the kernel
Strategy:	(Cooperatively) move data closer to SM: <ul style="list-style-type: none">• Shared Memory• (or Registers)• (or Constant Memory)• (or Texture Cache)



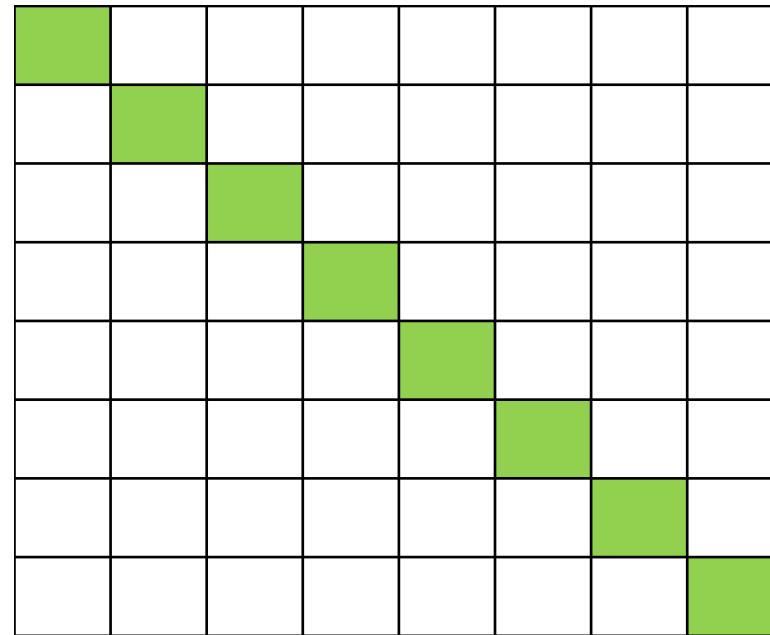
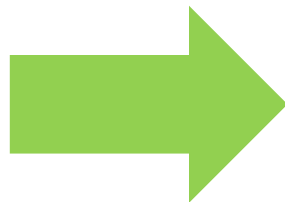
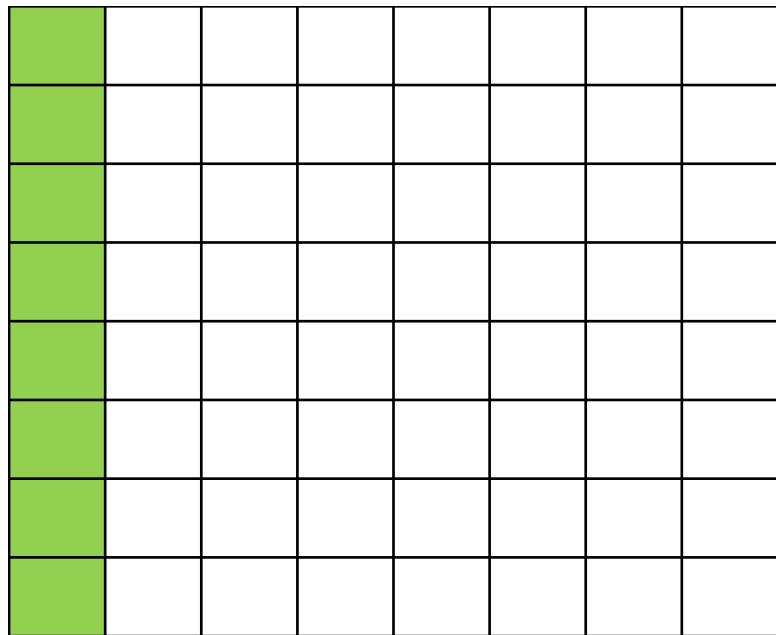
PERF-OPT QUICK REFERENCE CARD

Category:	Shared Mem Bandwidth Bound - Shared Memory
Problem:	Shared memory bandwidth bottleneck
Goal:	<u>Reduce</u> amount of data traffic to/from global mem
Indicators:	Shared memory loads or stores saturate
Strategy:	Reduce Bank Conflicts (insert padding) Move data from shared memory into registers Change data layout in shared memory



USING SHARED MEMORY WITHOUT CONFLICTS

FIXING BANK CONFLICTS



FIXING BANK CONFLICTS

```
// No bank-conflict transpose
// Same as transposeCoalesced except the first tile dimension is padded
// to avoid shared memory bank conflicts.
__global__ void transposeNoBankConflicts(float *odata, const float *idata)
{
    __shared__ float tile[TILE_DIM][TILE_DIM+1];

    int x = blockIdx.x * TILE_DIM + threadIdx.x;
    int y = blockIdx.y * TILE_DIM + threadIdx.y;
    int width = gridDim.x * TILE_DIM;

    for (int j = 0; j < TILE_DIM; j += BLOCK_ROWS)
        tile[threadIdx.y+j][threadIdx.x] = idata[(y+j)*width + x];

    __syncthreads();

    x = blockIdx.y * TILE_DIM + threadIdx.x; // transpose block offset
    y = blockIdx.x * TILE_DIM + threadIdx.y;

    for (int j = 0; j < TILE_DIM; j += BLOCK_ROWS)
        odata[(y+j)*width + x] = tile[threadIdx.x][threadIdx.y + j];
}
```

```
Device : Tesla P100-PCIE-16GB
Matrix size: 1024 1024, Block size: 32 8, Tile size: 32 32
dimGrid: 32 32 1. dimBlock: 32 8 1
```

Routine	Bandwidth (GB/s)
copy	442.94
shared memory copy	472.60
naive transpose	111.13
coalesced transpose	363.88
conflict-free transpose	465.00

```
umehta@nid02222:~/summerschool>
```

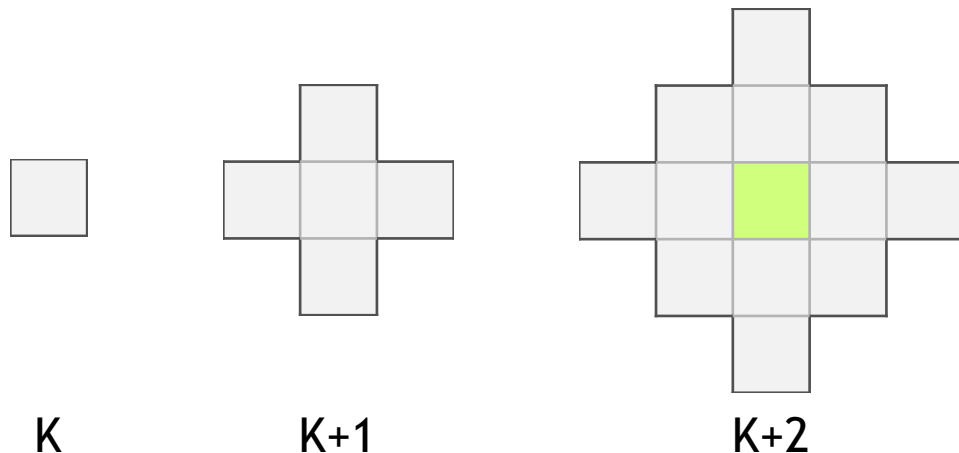
**ITERATION 3:
KERNELS WITH INCREASED
ARITHMETIC INTENSITY**

STENCIL OPERATIONS

2nd order vs 1st order

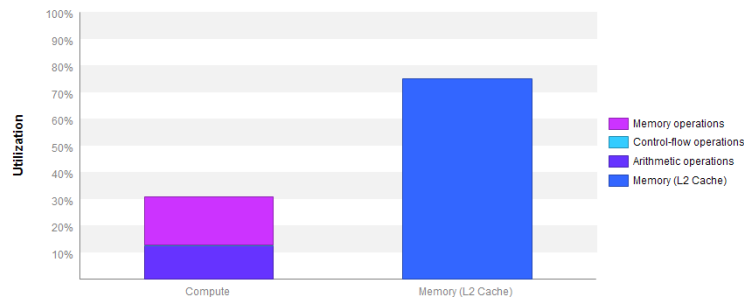
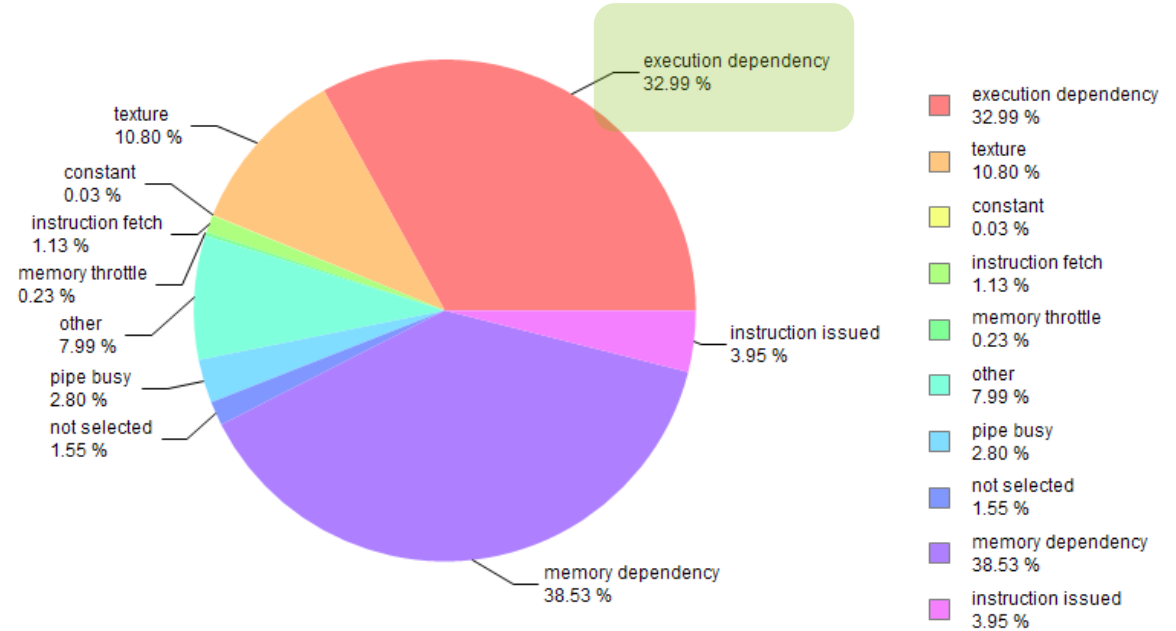
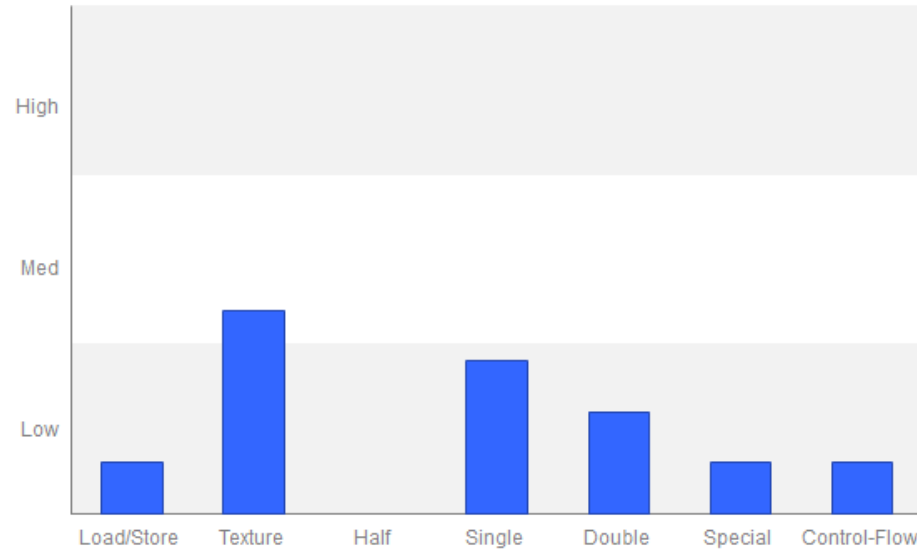
Performs 4x the FP operations

DRAM memory footprint is the same (assuming no over fetch)



FUNCTION UNIT UTILIZATION AND STALL REASONS

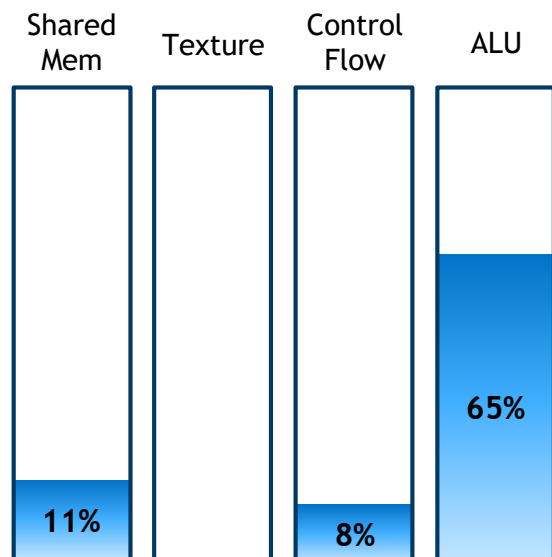
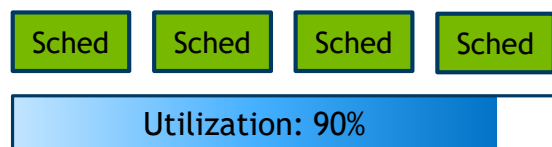
Execution Dependencies
starting to become
significant!



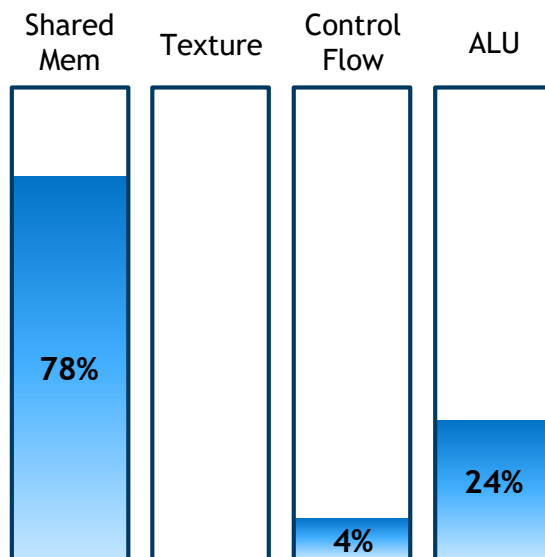
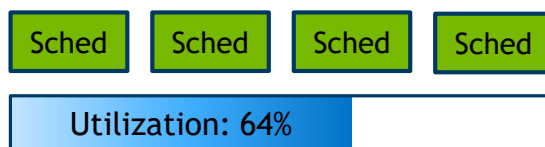
Functional units are not the bottlenecks in
Stencils, even with higher order stencils!

INSTRUCTION THROUGHPUT

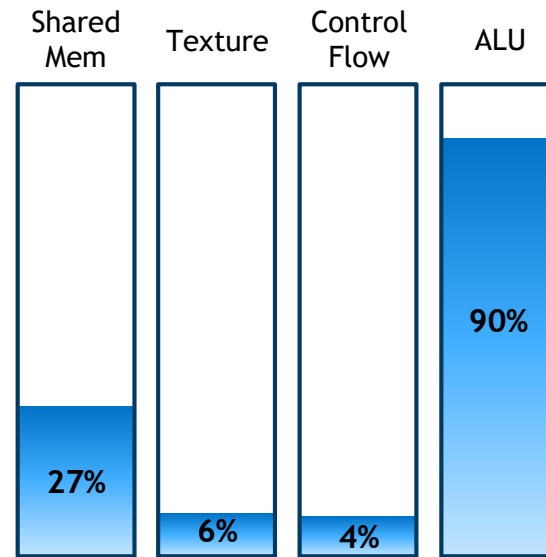
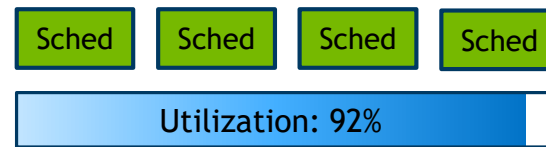
Schedulers saturated



FU saturated




Schedulers and FU saturated



STALL REASONS: EXECUTION DEPENDENCY

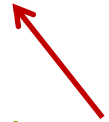
```
a = b + c; // ADD
```

```
d = a + e; // ADD
```



```
a = b[i]; // LOAD
```

```
d = a + e; // ADD
```



Memory accesses may influence execution dependencies

Global accesses create longer dependencies than shared accesses

Read-only/texture dependencies are counted in Texture

Instruction level parallelism can reduce dependencies

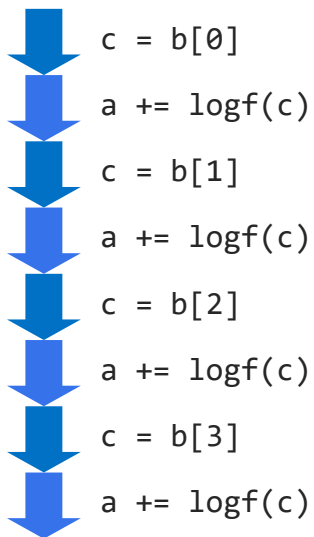
```
a = b + c; // Independent ADDs
```

```
d = e + f;
```

ILP AND MEMORY ACCESSSES

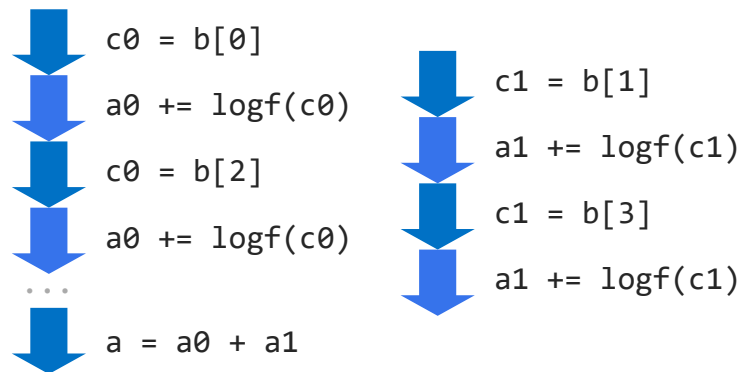
No ILP

```
float a = 0.0f;  
for( int i = 0 ; i < N ; ++i )  
    a += logf(b[i]);
```



2-way ILP (with loop unrolling)

```
float a, a0 = 0.0f, a1 = 0.0f;  
for( int i = 0 ; i < N ; i += 2 )  
{  
    a0 += logf(b[i]);  
    a1 += logf(b[i+1]);  
}  
a = a0 + a1
```



#pragma unroll is useful to extract ILP

Manually rewrite code if not a simple loop

PERF-OPT QUICK REFERENCE CARD

Category:	Bandwidth Bound - Register Caching
Problem:	Data is reused within threads and memory bw utilization is high
Goal:	<u>Reduce</u> amount of data traffic to/from global mem
Indicators:	High device memory usage, latency exposed Data reuse within threads and small-ish working set Low arithmetic intensity of the kernel
Strategy:	<ul style="list-style-type: none">• Assign registers to cache data• Avoid storing and reloading data (possibly by assigning work to threads differently)• Avoid register spilling



PERF-OPT QUICK REFERENCE CARD

Category:	Latency Bound - Instruction Level Parallelism
Problem:	Not enough independent work per thread
Goal:	Do more parallel work inside single threads
Indicators:	High execution dependency, increasing occupancy has no/little positive effect, still registers available
Strategy:	<ul style="list-style-type: none">• Unroll loops (#pragma unroll)• Refactor threads to compute n output values at the same time (code duplication)



PERF-OPT QUICK REFERENCE CARD

Category:	Compute Bound - Algorithmic Changes
Problem:	GPU is computing as fast as possible
Goal:	Reduce computation if possible
Indicators:	Clearly compute bound problem, speedup only with less computation
Strategy:	<ul style="list-style-type: none">• Pre-compute or store (intermediate) results• Trade memory for compute time• Use a computationally less expensive algorithm• Possibly: run with low occupancy and high ILP

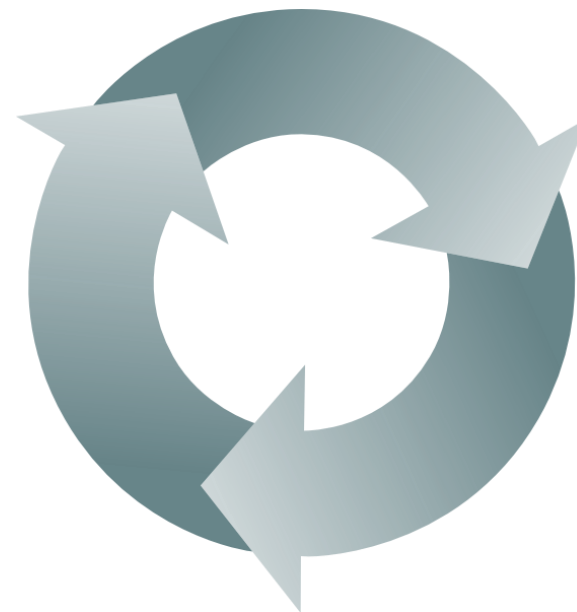


SUMMARY

SUMMARY

Performance Optimization is a Constant Learning Process

1. Know your application
2. Know your hardware
3. Know your tools
4. Know your process
 - Identify the Hotspot
 - Classify the Performance Limiter
 - Look for indicators
5. Make it so!



REFERENCES

CUDA Documentation

Best Practices: <http://docs.nvidia.com/cuda/cuda-c-best-practices-guide/>

Kepler Tuning Guide: <http://docs.nvidia.com/cuda/kepler-tuning-guide>

Maxwell Tuning Guide: <http://docs.nvidia.com/cuda/maxwell-tuning-guide>

Pascal Tuning Guide: <http://docs.nvidia.com/cuda/pascal-tuning-guide>

Parallel ForAll devblog

<http://devblogs.nvidia.com/parallelforall/>

GTC Sessions:

<http://on-demand-gtc.gputechconf.com>

THANK YOU

