



WC: A New GPU Programming Model

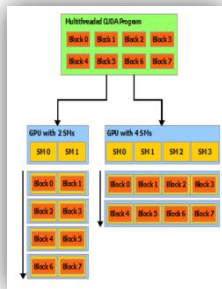
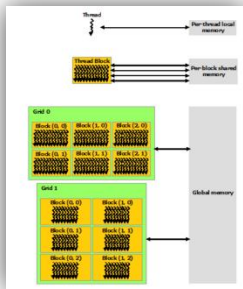
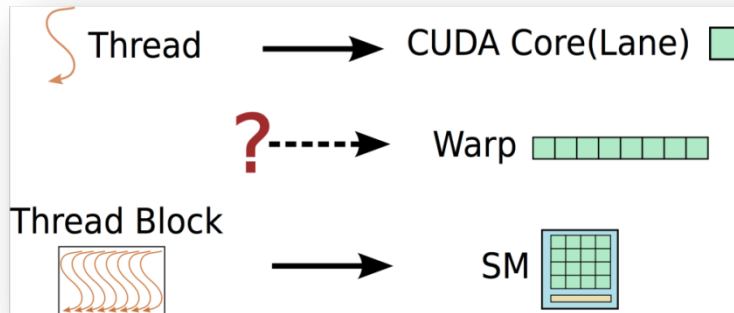
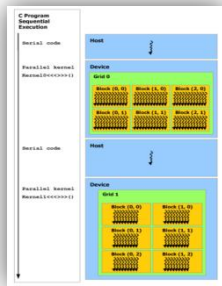
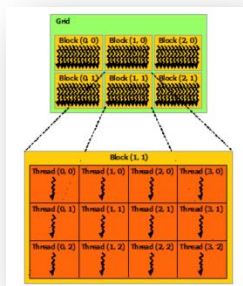
Ang Li, Pacific Northwest National Laboratory (PNNL)

Weifeng Liu (Norwegian University of Science and Technology)

Linnan Wang (Brown University)

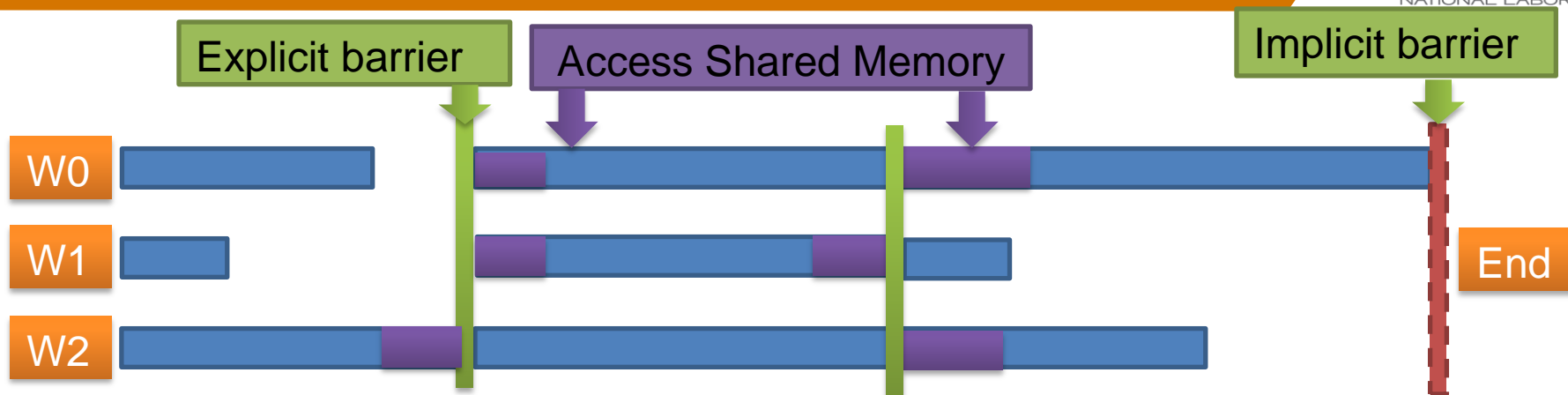
Kevin Barker and Shuaiwen Leon Song (PNNL)

The Missing Warp in CUDA

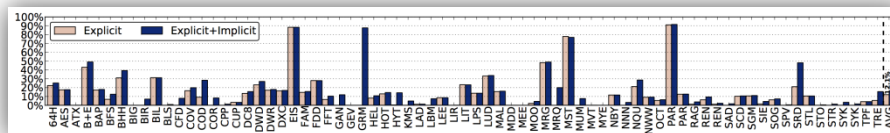


- ▶ CUDA's 2-level programming model
 - Thread-block → SM, Thread → CUDA core
- ▶ GPU Hardware's 3-level execution model
 - SM, warp, SIMD-lane
- ▶ CUDA programming/execution model mismatch
- ▶ Warp cannot be simply ignored
 - Warp divergence
 - Inter-warp synchronization

Inter-warp Synchronization



- ▶ Explicit Synchronization
- ▶ Implicit Synchronization



12% (up to 90%) explicit stalls, 4.4% implicit stalls for 80 GPGPU applications on P100

Large & correlated design space

- ▶ Difficult to program
 - Setting kernel configuration `<<<grid(x,y,z), block(x,y,z)>>>`
 - Related to resources allocation: register, shared-mem, warp-slots
- ▶ Difficult to optimize
 - Without knowing “warp” it’s hard to achieve optimal performance
 - Thinking about warp, program about warp, but never express warp
- ▶ Synchronization overhead
 - Warp divergence/unbalancing
 - Explicit warp synchronization
 - Implicit warp synchronization

CUDA’s programming/execution model mismatch

How to resolve the “mismatch” issue?

- ▶ Warp-Consolidation Model
 - 5 advantages over CUDA model
 - 3 code examples
 - 2 drawbacks & solution
 - 2 comparisons with other techniques
- ▶ Evaluation
 - Settings
 - Results
- ▶ Potential Extension
 - Active context switch
 - Volta



▶ WC Model

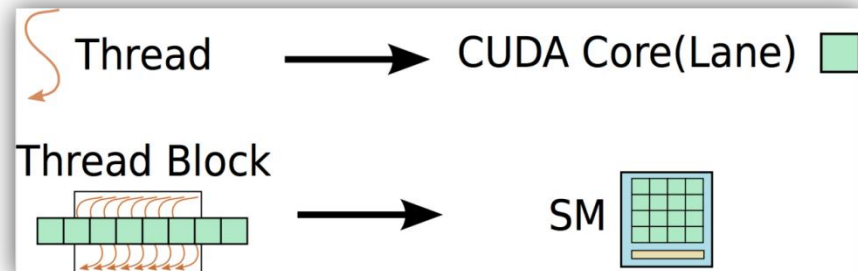
- Unify TB and warp: 1-warp/thread-block
- 2-level programming/execution model

▶ Combine the advantages

- thread-block(free resource allocation)
- warp (fast communication, synchronization and cooperation)

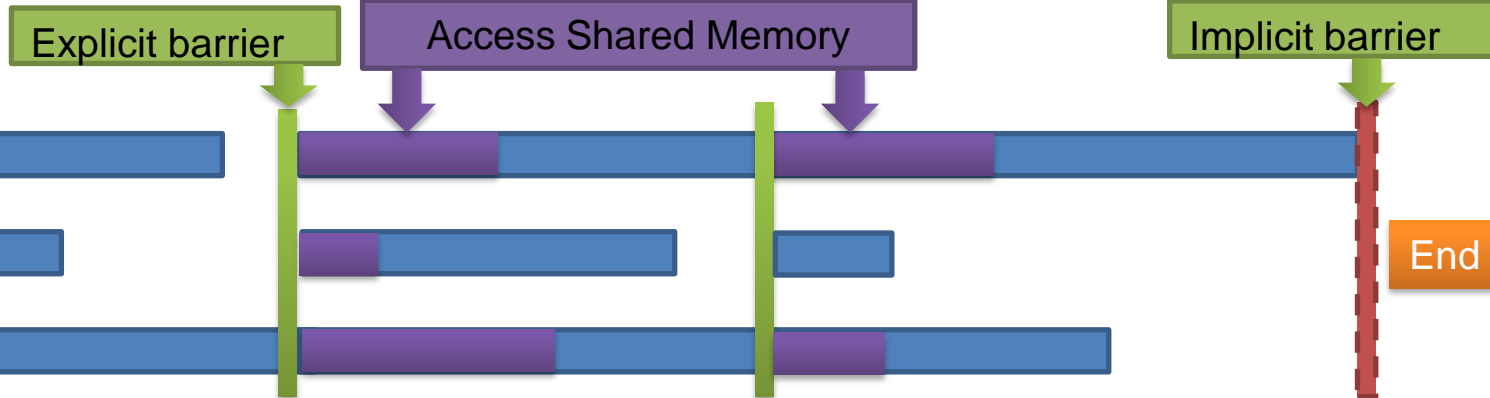
▶ Motivates novel SMT-like programming concept for GPU

- Partition computation space into multiple independent warp based jobs. A job is handled by a warp
- Threads in a TB/warp can efficiently sync, cooperate and communicate

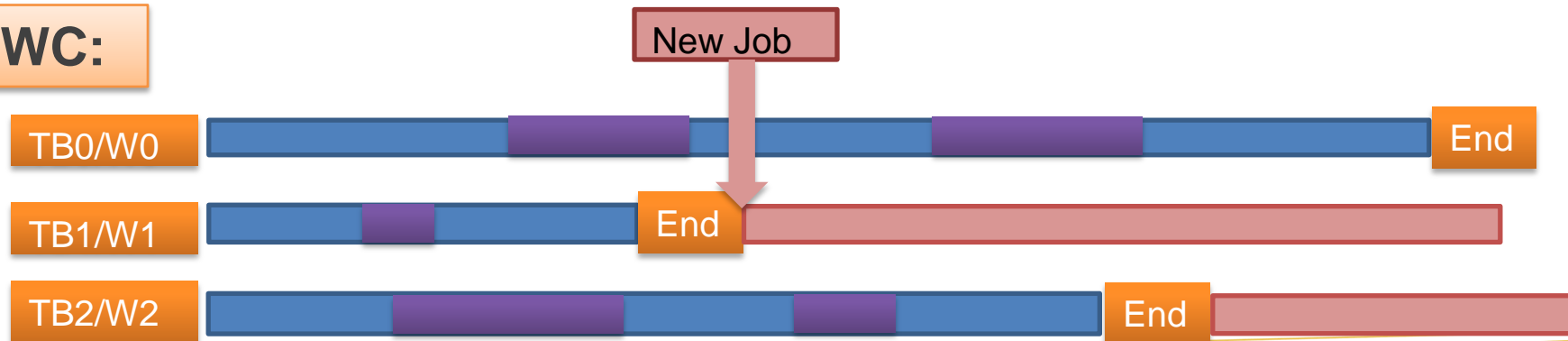


Adv 1: No synchronization & barrier

CUDA:



WC:



▶ *CUDA Programming Model*

Hard to program

- *Large design space <<<grid(x,y,z), block(x,y,z)>>>*
- *Correlated resources usage: register, shared-mem, warp-slots, etc.*

▶ WC Programming Model

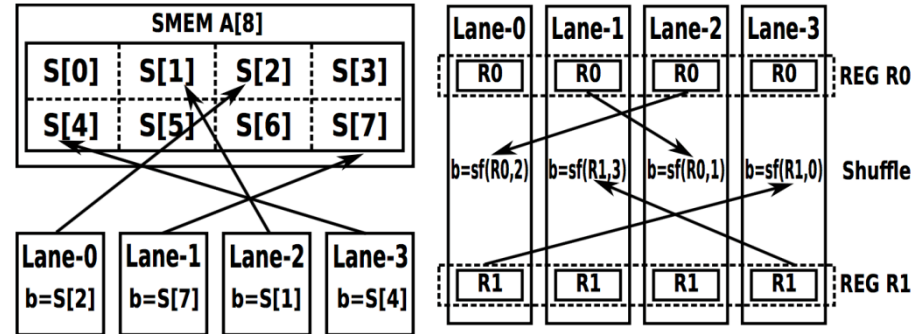
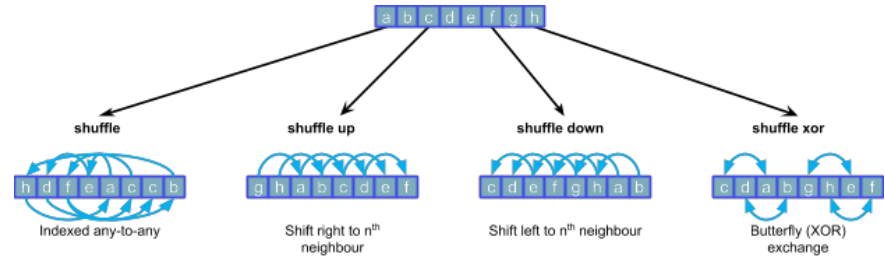
Easy to program

- Fixed block configuration <<<grid(x,y,z), 32>>>
- Fixed per-warp resources → No correlated resource allocation
- Significantly reduces design space
- CPU-like programming: simultaneous multi-threaded SSE code

Adv 3: Register fast communication

- ▶ Shared memory communication:
 - reg → write → sync → read → reg
- ▶ Register shuffle communication:
 - reg → shuffle → reg
 - Very flexible
- ▶ Communication
 - CUDA: collective shared memory communication
 - WC: distributive message passing communication

Flexible Shuffle

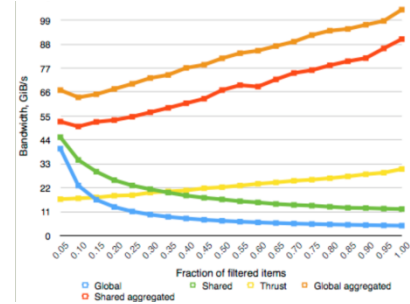


Adv 4: Fine-grained cooperation

`__any()`
`__all()`
`__ballot()`
`__ffs()`
`__clz()`
`__popc()`

Warp-aggregated Atomics

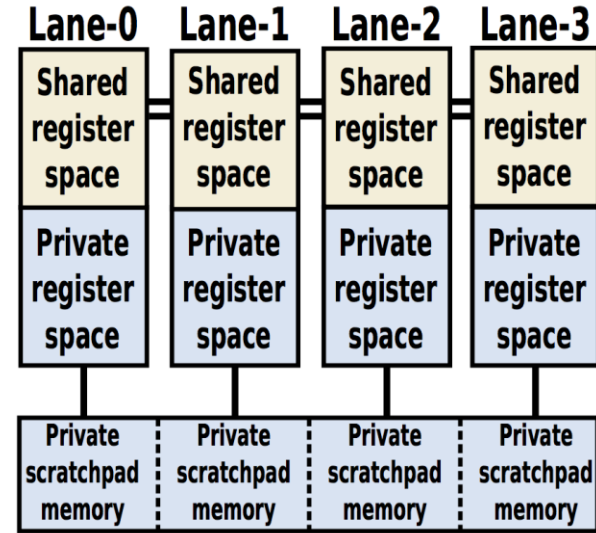
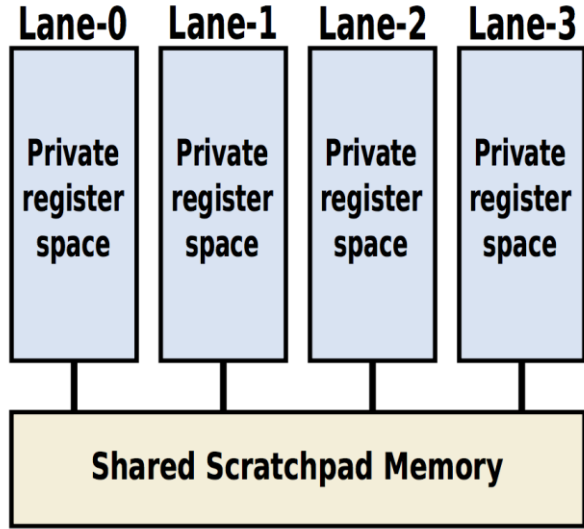
```
1 __device__
2 int atomicAggInc(int *ctr) {
3     int mask = __ballot(1);
4     int leader = __ffs(mask) - 1; // select the leader
5     int res;
6     if(lane_id() == leader) // leader does the update
7         res = atomicAdd(ctr, __popc(mask));
8     res = __shfl(res, leader);
9     // each thread computes its own value
10    return res + __popc(mask & ((1 << lane_id()) - 1));
11 }
```



[1] M. Bauer et al, *Singe: Leveraging warp specialization for high performance on GPUs*, PPOPP-14

[2] Andy Adinets, *Optimizing Filtering with Warp-Aggregated Atomics*, <https://devblogs.nvidia.com/cuda-pro-tip-optimized-filtering-warp-aggregated-atomics/>

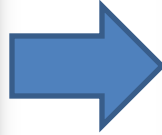
Adv 5: Extended Register Space



- ▶ Register spilling in shared memory
 - Register shuffling relax the usage of shared memory as comm buffer
 - Shared memory allocation is now private a warp
 - No shared memory bank conflict

Exp 1: How to program in WC

```
__global__ void Kernel(...){
  __shared__ sh[N];
  for(...){
    __syncthreads();
    __threadfence_block();
    ...}
}
num_thds = 512;//16 warps
num_blks = (int)ceil(num_nodes/
             (double)num_thds;
dim3 grid(num_blks,1,1);//grid config
dim3 thds(num_thds,1,1);//blk config
do{
  Kernel<<<grid,thds>>>(...)
while(stop);
```



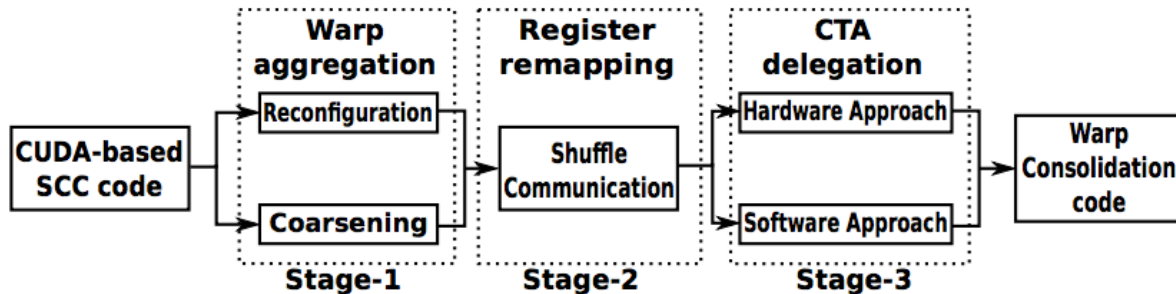
```
__global__ void Kernel(...){
  __shared__ volatile sh[N];//volatile smem
  for(...){
    //sync is avoided
    //blockwise memory fence is avoided
    ...}
}
num_thds = 32;//reduce to one warp
num_blks = (int)ceil(num_nodes/
             (double)num_thds;//num of blks x16
dim3 grid(num_blks,1,1);//grid config
dim3 thds(num_thds,1,1);//blk config
do{
  Kernel<<<grid,thds>>>(...)
while(stop);
```

- Synchronization is removed
- Shared memory allocate as volatile
- Thread block configuration == 32

Exp 2: How to transform legacy code

```
if(i0<n_vectors) //replicate or use "for"
  d_output[i0]=(float)X*k_2powneg32;
unsigned v_log2stridem1=v[__ffs(stride)-2];
unsigned v_stridemask=stride-1;
//"for" statement need to replicated,
//can be fused later if possible
for(int i=i0+stride;i<n_vectors;i+=stride){
  X^=v_log2stridem1^v[__ffs(~((i-stride)|
    v_stridemask))-1];
  d_output[i]=(float)X*k_2powneg32;
}}
```

```
if(i0_0<n_vectors)
  d_output[i0_0]=(float)X_0*k_2powneg32;
if(i0_1<n_vectors)
  d_output[i0_1]=(float)X_1*k_2powneg32;
unsigned v_log2stridem1=v[__ffs(stride)-2];
unsigned v_stridemask=stride-1;
for(int i=i0_0+stride;i<n_vectors;i+=stride){
  X_0^=v_log2stridem1^v[__ffs(~((i-stride)|
    v_stridemask))-1];
  d_output[i]=(float)X_0*k_2powneg32;
}
for(int i=i1_0+stride;i<n_vectors;i+=stride){
  X_1^=v_log2stridem1^v[__ffs(~((i-stride)|
    v_stridemask))-1];
  d_output[i+32]=(float)X_1*k_2powneg32;
}}
```



Input : Original Kernel Code

Output : Warp-Consolidation based Optimized Kernel Code

Stage-1 Warp Aggregation(kernel code):

```
if CTA-size can be adjusted to 32 threads then
  Adjust CTA-size to 32 threads;
else
  Perform aggressive warp coarsening;
  for each code block separated by __syncthreads() do
    Replicate each var reused across the code block;
    Replicate statements depending on threadIdx.x;
    Convert "if" statements to "for" statements;
    Replicate "for" and "while" block;
    Perform loop fusion if possible;
  end
end
```

end

Remove all sync statements;

return optimized code after warp-aggregation;

Stage-2 Register Remapping(kernel code):

```
if Fixed s-mem access pattern and no addr calculation then
  Remove all shared memory allocation;
  Partition smem space into warp-based chunks;
  Allocate these chunks in different registers;
  Replace smem references by shuffle instructions;
end
```

end

return optimized code after register-remapping;

Stage-3 Warp Delegation(kernel code):

```
Include the Warp.Delegation.cuh header file;
Pack kernel body into the WARP_DELEGATION macro;
Set partition strategy;
Convert kernel invocation;
return optimized code after warp-delegation;
```

Exp 3: How to communicate via reg

```
for(int k=1;k<nz-1;k++){
float a_left_right,a_up,a_down;
if((i<nx)&&(j<ny))
    top=A0[Index3D(nx,ny,i,j,k+1)];
if(w_region){
    ...
    a_left_right=x_l_bound?A0[Index3D(
        nx,ny,i-1,j,k)]:sh_A0[sh_id-1];
    Anext[Index3D(nx,ny,i,j,k)]=(top+bottom
        +a_up+a_down+sh_A0[sh_id+1]
        +a_left_right)*c1-sh_A0[sh_id]*c0;
}
```

```
__syncthreads();
bottom=sh_A0[sh_id];
sh_A0[sh_id]=top;
__syncthreads();
```

```
//shuffle is outside branch
a_left=__shfl_up(r0_A0,1);//left neighbor
a_right=__shfl_down(r0_A0,1);//right neighbor
a_other=__shfl(r1_A0,0);//2nd reg space
if(w_region){
    ...
    a_left=x_l_bound?A0[Index3D(nx,ny,i-1,j,k)]
        :a_left;
    a_right=(threadIdx.x==31)?a_other:a_right;
    Anext[Index3D(nx,ny,i,j,k)]=(top+bottom
        +a_up+a_down+a_left
        +a_right)*c1-r0_A0*c0;
}
...
bottom=r0_A0;
r0_A0=top;
```

- Shuffle is very flexible (up, down, shfl)
- Be careful when shuffle with warp divergence
- No synchronization required

Drawback 1: Occupancy degradation

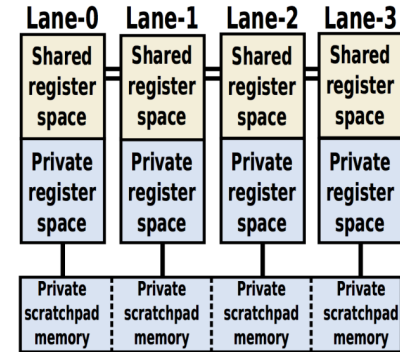
GPU	Architecture	SMs	Thread-Blocks/SM	Warps/SM
Tesla-K80	Kepler	15	16	64
Tesla-M40	Maxwell	24	32	64
Tesla-P100	Pascal	56	32	64
Tesla-V100	Volta	80	32	64

- ▶ Occupancy is at maximum 0.5
 - Resolved by warp-delegation (a variant of [1])
 - CUDA Best Practice Guide (50% occupancy is sufficient)
 - We hope $\text{TBs/SM} == \text{Warps/SM}$ in future GPUs

[1] A. Li et al. "Locality-Aware CTA Clustering for Modern GPUs." in *ASPLOS-17, ACM*

Drawback 2: Shared Mem/Reg usage

- ▶ Extra shared mem usage
 - Ideally can be mitigated & released by register shuffling
 - Otherwise, multiplexing can be an alternative solution[1]
- ▶ Extra register usage
 - Our experiment show that reg usage not increase much
 - Can be resolved by shared memory spilling



[1] Y. Yang et al. "Shared memory multiplexing: a novel way to improve GPGPU throughput" in *PACT-12, ACM*

Comparison 1: Thread Coarsening

- ▶ Fuse multi-threads so per-thread workload increases but thread number decrease

Thread Coarsening

Increase ILP

Register reuse

Reduce auxiliary inst

Coarsening factor

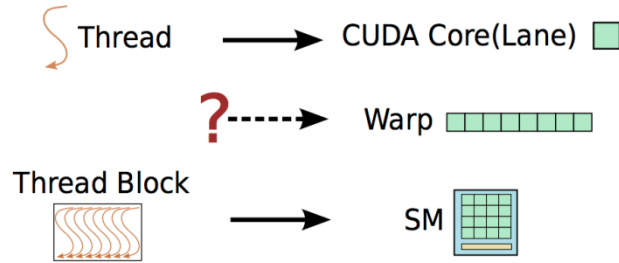
Warp Consolidation

An approach to reduce warps/TB to 1 for legacy CUDA code

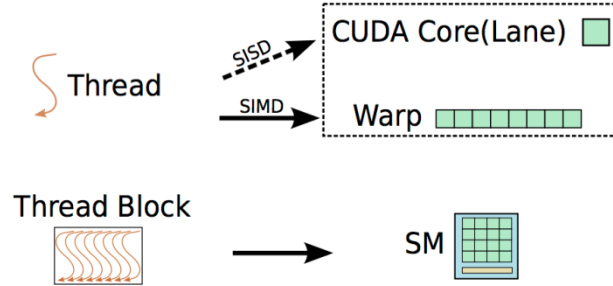
Coarsening === 1

Comparison 2: Warp Specialization

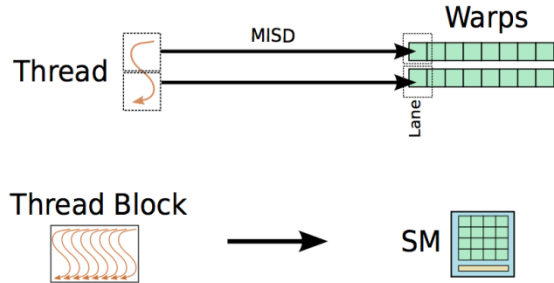
CUDA:



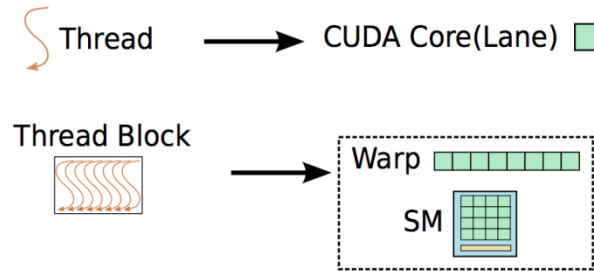
Warp-Centric [2]:



Warp-Specialization [1]:



Warp Consolidation (this work):



[1] M. Bauer et al. "Singe: leveraging warp specialization for high performance on GPUs." in *PPoPP-14, ACM*

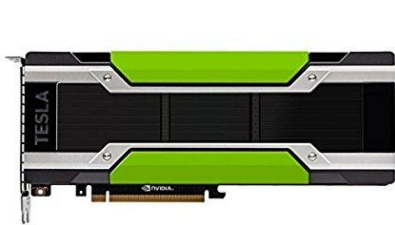
[2] S. Hong et al. "Accelerating CUDA graph algorithms at maximum warp" in *PPoPP-11, ACM*.

How to resolve the “mismatch” issue?

- ▶ Warp-Consolidation Model
 - 5 advantages over CUDA
 - 3 examples
 - 2 drawbacks & solution
 - 2 comparisons with other techniques
- ▶ Evaluation
 - Settings
 - Results
- ▶ Potential Extension
 - Active context switch
 - Volta



Evaluation Settings



**Tesla-K80
Kepler**



**Tesla-M40
Maxwell**



**Tesla-P100
Pascal**

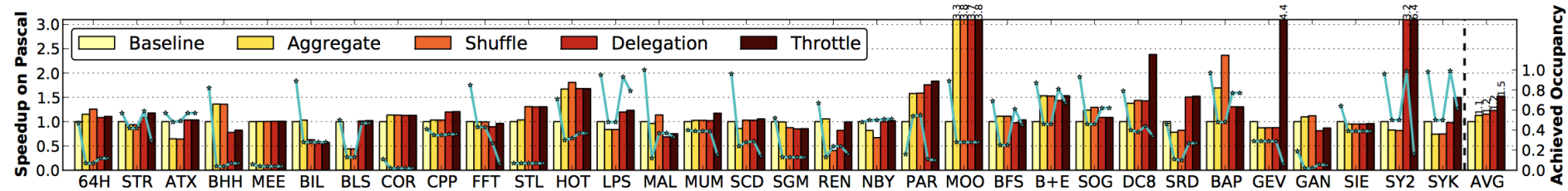


**Tesla-V100
Volta**

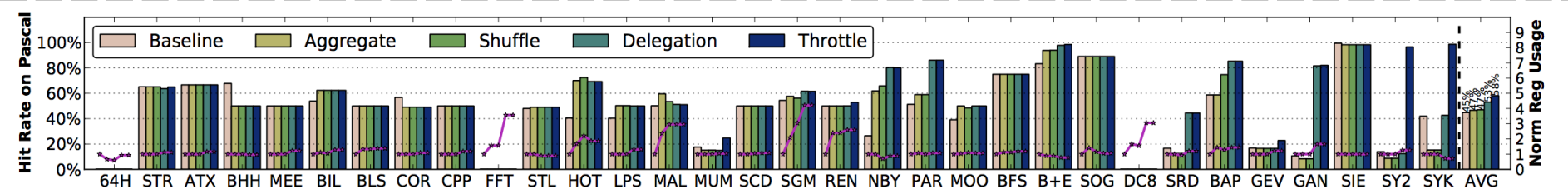
Application	Description	abbr.	Kernel Name	CTAs	WPs	Regs	SMem	Sync	Comm	Ref
<i>g4h</i>	64 and 256 bit histogram calculation	G4H	<i>mergeHistogram2dKernel()</i>	256	8	15	1000B	Y	N	[32]
<i>streamcluster</i>	Assigning points to nearest centers	STR	<i>kernel_compute_cost()</i>	128	16	25	0B	N	N	[33]
<i>atax</i>	Matrix transpose and vector multiply	ATX	<i>atax_kernel2()</i>	256	8	18	0B	N	N	[34]
<i>lh</i>	Gravitational forces via Barnes-Hut method	BHR	<i>BoundingBoxKernel()</i>	168	32	32	24000B	Y	Y	[35]
<i>mersenne</i>	Mersenne Twister random generator	ME	<i>BlockMulKernelGPU()</i>	32	4	18	0B	N	N	[32]
<i>binomial</i>	Option call price via binomial model	BTL	<i>binomialOptionsKernel()</i>	1024	8	17	2007B	Y	Y	[32]
<i>BlackScholes</i>	Option call price via Black-Scholes model	BLF	<i>BlackScholesGPU()</i>	480	4	21	0B	N	N	[32]
<i>corr</i>	Correlation computation	CR	<i>corr_kernel()</i>	8	8	22	0B	N	N	[34]
<i>cutcp</i>	Compute Coulombic potential for 3D grid	CPP	<i>cutcp_kernel()</i>	512	4	22	0B	N	N	[36]
<i>fft</i>	Fast Fourier transform	FFT	<i>FFT32_device()</i>	32768	2	56	4500B	Y	Y	[37]
<i>stencil</i>	Jacobi stencil operation on regular 3D grid	STL	<i>block2D_hybrid_cudaconv.at()</i>	64	4	29	1000B	Y	Y	[36]
<i>hotspot</i>	Estimate processor temperature	HTP	<i>calculate_temp()</i>	7596	8	38	3000B	Y	Y	[33]
<i>lps</i>	3D Laplace solver	LPS	<i>GPU_Laplace3d()</i>	5000	4	16	2300B	Y	Y	[38]
<i>matrixMul</i>	Matrix multiplication	MAL	<i>matrixMulCUDA()</i>	16384	32	27	8000B	Y	Y	[32]
<i>mem</i>	Pair-wise local sequence alignment for DNA	MEM	<i>memalignGPUKernel()</i>	196	8	23	0B	N	N	[38]
<i>scalarProd</i>	Scalar products of input vector pairs	BCD	<i>scalarProdGPU()</i>	4096	8	24	4000B	Y	Y	[32]
<i>sgemm</i>	Single precision general matrix multiply	SGM	<i>mysgemvNT()</i>	496	4	46	512B	Y	Y	[36]
<i>recursiveGaussian</i>	Recursive Gaussian filter	RGR	<i>d_rangepool()</i>	1024	8	10	1062B	Y	Y	[32]
<i>nbody</i>	All-pairs gravitational N-body simulation	NBT	<i>integrateBodySet()</i>	224	8	36	4000B	Y	Y	[32]
<i>pathfinder</i>	Dynamically finding a path in 2D grid	PAF	<i>dynproc_kernel()</i>	3334	1	13	256B	Y	Y	[33]
<i>MonteCarlo</i>	Option call price via Monte-Carlo method	MOC	<i>MonteCarloReplacer()</i>	1024	8	29	4000B	Y	Y	[32]
<i>lfs</i>	Breadth first search	BFS	<i>Kernel()</i>	1954	16	16	0B	Y	N	[33]
<i>B+tree</i>	B+tree Operation	B+E	<i>findKi()</i>	10000	8	27	0B	Y	N	[33]
<i>SobelCGRNG</i>	Sobel edge detection filter for images	BOG	<i>sobelGPU_kernel()</i>	25600	2	19	128B	Y	Y	[32]
<i>dstfft</i>	Discrete cosine transform for 8-bit block	DCT	<i>CUDAkernel(DCT)</i>	4096	2	18	500B	Y	Y	[32]
<i>rsad</i>	Speckle reducing anisotropic diffusion	SRD	<i>reducer()</i>	2048	16	25	4000B	Y	Y	[33]
<i>backprop</i>	Perceptron back propagation	BAP	<i>lapon_layerforward_CUDA()</i>	32768	8	18	1062B	Y	Y	[33]
<i>gemmm</i>	Scalar vector and matrix multiplication	GEV	<i>gemmm_kernel()</i>	128	8	23	0B	N	N	[34]
<i>gaussian</i>	Solving variables in a linear system	G4R	<i>Fun()</i>	2	16	12	0B	N	N	[33]
<i>single</i>	Monte Carlo single Asian option	BIE	<i>initRVG()</i>	782	4	32	0B	Y	Y	[32]
<i>sv2k</i>	Symmetric rank-2 operations	SV2	<i>sv2k_kernel()</i>	16384	8	19	0B	N	N	[34]
<i>svk</i>	Symmetric rank-k operations	SVK	<i>svk_kernel()</i>	4096	8	28	0B	N	N	[34]

32 Applications from commonly used GPGPU benchmark suits

Evaluation Results



On average 1.7x, 2.3x, 1.5x and 1.2x over Tesla K80 (Kepler), M40 (Maxwell), P100 (Pascal) and V100 (Volta)



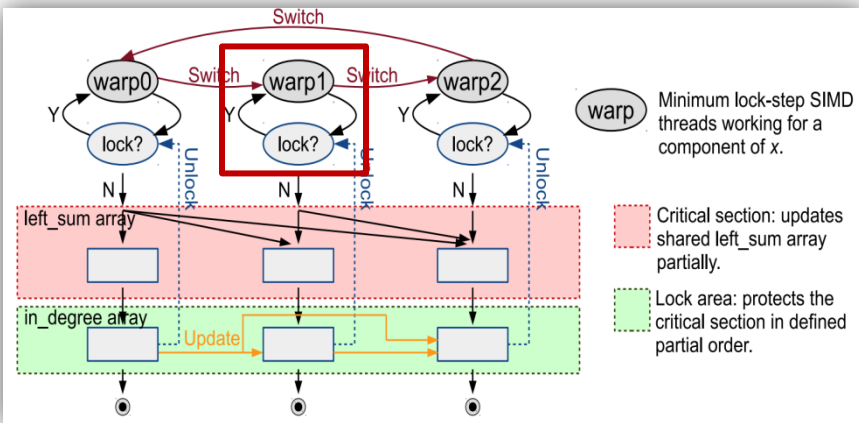
Performance gain is not directly from cache by reducing sync, communication and cooperation (SCC) overhead

How to resolve the “mismatch” issue?

- ▶ Warp-Consolidation Model
 - 5 advantages over CUDA
 - 3 examples
 - 2 drawbacks & solution
 - 2 comparisons with other techniques
- ▶ Evaluation
 - Settings
 - Results
- ▶ Potential Extension
 - Active context switch
 - Volta



Extension 1: Active context switch



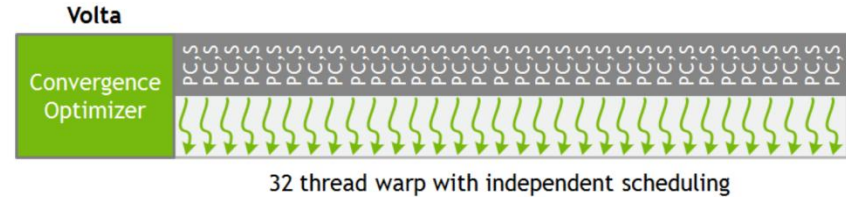
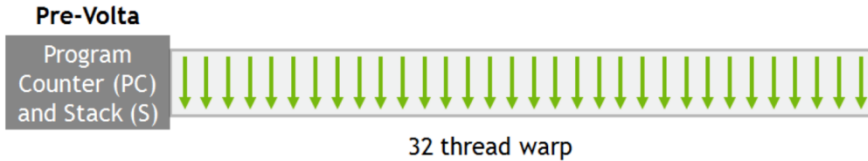
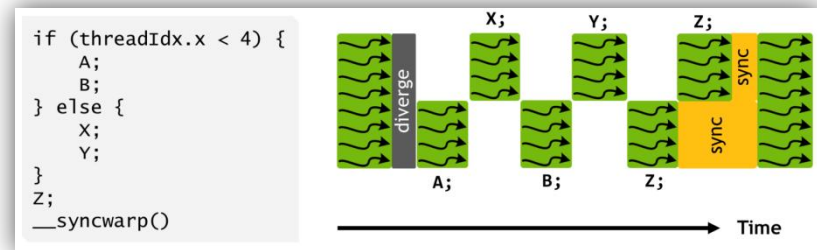
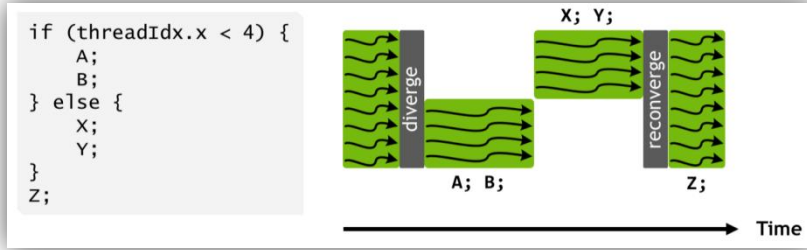
```
while s_in_degree[i] + 1  $\neq$  d_in_degree[i] do  
    //busy wait  
end while
```

Not a deadlock if contains
“clock()” which appear to
signal warp context switch

- Exploit intra-TB and inter-TB data reuse
 - Invoke clock() after memory fetch when there is more inter-TB locality

Weifeng Liu et al. **Fast synchronization-free algorithms for parallel sparse triangular solves with multiple right-hand sides.** *Concurrency and Computation: Practice and Experience*, 2017.

Extension 2: Volta opportunities



- ▶ Threads in a warp can proceed in sub-warp granularity
 - Benefit: more flexible warp sync, communication and cooperation
 - Issue: threads in a warp has to synchronize

Extension 2: Volta opportunities



```
// Calculate AB with NVIDIA Tensor Cores
// Kernel executed by 1 Warp (32 Threads)
__global__ void tensorOp(float *D, half *A, half *B) {
    // 1. Declare the fragments
    wmma::fragment<wmma::matrix_a, M, N, K, half, wmma::col_major> Amat;
    wmma::fragment<wmma::matrix_b, M, N, K, half, wmma::col_major> Bmat;
    wmma::fragment<wmma::accumulator, M, N, K, float, void> Cmat;
    // 2. Initialize the output to zero
    wmma::fill_fragment(Cmat, 0.0f);
    // 3. Load the inputs into the fragments
    wmma::load_matrix_sync(Amat, A, M);
    wmma::load_matrix_sync(Bmat, B, K);
    // 4. Perform the matrix multiplication
    wmma::mma_sync(Cmat, Amat, Bmat, Cmat);
    // 5. Store the result from fragment to global
    wmma::store_matrix_sync(D, Cmat, M, wmma::mem_col_major);
}
```

How can the WC model behaves to be more efficiently operates the Tensor Core!

- ▶ **Warp-Consolidation:** a GPU Programming and Execution model that
 - Unifies warp and thread block (no explicit & implicit sync)
 - Communicates via register while cooperates via warp voting
- ▶ **Applicability:**
 - Simplified programming model than CUDA
 - SCC (sync, communication, cooperation) applications
- ▶ **1.7x, 2.3x, 1.5x and 1.2x** average speedups across 32 GPGPU applications on Kepler, Maxwell, Pascal and Volta GPUs, respectively



U.S. DEPARTMENT OF
ENERGY



WC: A New GPGPU Programming Model

Ang Li, Pacific Northwest National Laboratory (PNNL)

Weifeng Liu (Norwegian University of Science and Technology)

Linnan Wang (Brown University)

Kevin Barker and Shuaiwen Leon Song (PNNL)



Pacific Northwest
NATIONAL LABORATORY

Proudly Operated by Battelle Since 1965