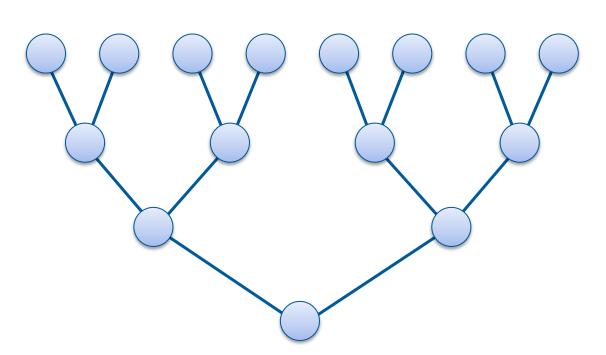


Parallel Reduction



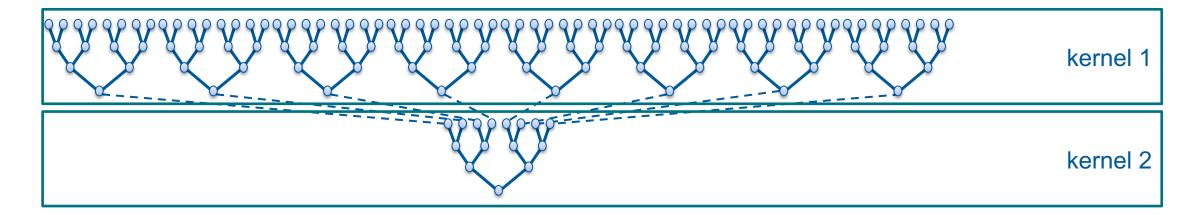
- Common and important data parallel primitive
 - Examples: minimum, maximum, sum
- Many data elements → single output (associative!)
- Easy to implement in CUDA
- Tree-based approach



Parallel Reduction



- Need to use multiple thread blocks
 - To process large arrays
 - To fully utilize GPU
- Partition the array, one block per partition
- How to communicate partial results?



Parallel Reduction #1: Atomic Global



```
<u>_global__</u>        void reduce_atomic_global(const float* input, float* result, int elements) {
  int id = threadIdx.x + blockIdx.x*blockDim.x;
  atomicAdd(result, input[id]);
```



	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic globa	635.355 ms	1.524 GB/s	1.000	1.000	
C	CPU baseline: 283.989 ms PU parallel baselin 85.393 ms	e:			

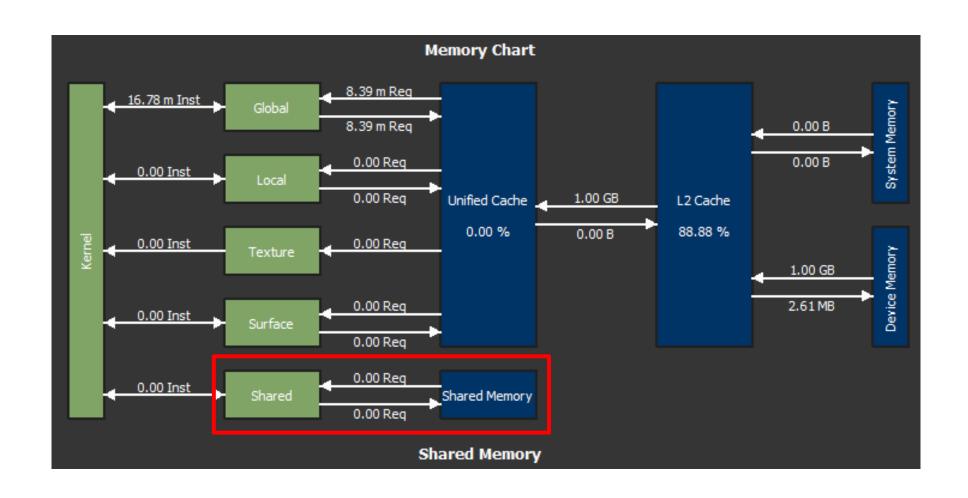
Parallel Reduction #1: Atomic Global



	Instructions	Requests	% Peak	Bank Conflicts						
Shared Load	instructions 0	nequests 0	76 PEAK 0	Dank Conflicts	0					
Shared Store	0	0	ō		0					
Shared Atomic	0									
Total	0	0	0		0					
10101									First Level Co.	
									First-Level Cac	
	Instructions	SM->TEX Requests	% Peak	Hit Rate	TEX->L2 Requests	% Peak	L2->TEX Returns	% Peak	TEX->SM Returns	% Peak
Global Load Cached	8,388,608	8,388,608	0.01							
Global Load Uncached							33,554,432	0.05	16,777,216	0.02
Local Load Cached										
Local Load Uncached										
Surface Load									0	0
Texture Load										
Global Store					0	0				
Local Store										
Confees Chara	^	^	^	^	^	^				
Global Reduction	8,388,608	8,388,608	0.01		268,435,456	0.02				
Juriace Neutron										
Global Atomic		0	0	0	0	0	0	0	see above	see above
Global Atomic Cas										
Surface Atomic		0	0	0	0	0	0	0	see above	see above
Surface Atomic Cas	0				<u> </u>	<u> </u>				
Loads	8,388,608	8,388,608	0.01				33,554,432	0.05	16,777,216	0.02
Stores										
Total	16,777,216	16,777,216	0.02		268,435,456	0.02	33,554,432	0.05	16,777,216	0.02
									Second-Level Ca	iche
	TEX->L2 Requests	% Peak L2-	> TEX Returns	% Peak	Total Bytes	Total T	hroughput			
Global Load Cached					1,073,741,824		1,519,729,736.71			
Global Load Uncached										
Local Load Cached			33,554,432	0.05						
Local Load Uncached										
Surface Load										
Texture Load										
Global Store										
Local Store										
C	^									
Global Reduction	268,435,456	0.02			8,589,934,592		12,157,837,893.66			
		0.02			0,505,155,155					
Surface Neduction		V			<u></u>		<u> </u>			
Global Atomic		v	-		• • • •		0			
	0	0			0		0			
Global Atomic		0			0					
Global Atomic Global Atomic Cas	0 0	v			• • • •		0			
Global Atomic Global Atomic Cas Surface Atomic		0			0					
Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas	0	0	0	0	0		0			
Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads	0	0	0	0.05	0 0 1,073,741,824		0 1,519,729,736.71			
Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores		0 0 - 0	33,554,432 -	0 0.05 -	0 0 1,073,741,824 0		0 1,519,729,736.71 0		Device Memor	v
Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores	0 - 0 268,435,456	0 0 - 0 0.02	33,554,432 - 33,554,432	0 0.05 - 0.05	0 0 1,073,741,824 0 9,663,676,416		0 1,519,729,736.71 0		Device Memor	γ
Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores Total	0 - 0 268,435,456 L2<->FB Sectors	0 0 - 0 0.02 % Peak	0 33,554,432 - 33,554,432 Bytes	0 0.05 - 0.05 Thro	0 0 1,073,741,824 0 9,663,676,416 ughput		0 1,519,729,736.71 0		Device Memor	y
Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores Total	0 - 0 268,435,456 L2<->FB Sectors 33,572,397	0 0 - 0 0.02 % Peak 0.23	33,554,432 - 33,554,432 Bytes 1,074,316,70	0 0.05 - 0.05 Thro	0 0 1,073,741,824 0 9,663,676,416 ughput 1,520,543,398.07		0 1,519,729,736.71 0		Device Memor	y
Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores Total	0 - 0 268,435,456 L2<->FB Sectors	0 0 - 0 0.02 % Peak	0 33,554,432 - 33,554,432 Bytes	0.05 - 0.05 Thro 4	0 0 1,073,741,824 0 9,663,676,416 ughput		0 1,519,729,736.71 0		Device Memor	у

Parallel Reduction #1: Atomic Global







```
<u>global</u> void reduce_atomic_shared(const float* input, float* result, int elements) {
  int id = threadIdx.x + blockIdx.x*blockDim.x;
  float in = input[id];
  __shared__ float x;
  x = 0.0f;
  __syncthreads();
  atomicAdd(&x, in);
  __syncthreads();
  if (threadIdx.x == 0)
      atomicAdd(result, x);
```



		Time (2 ²⁸ floats)	В	andwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic gl	obal	635.355 ms		1.524 GB/s	1.000	1.000	
atomic sha	ared	26.911 ms	3	85.764 GB/s	23.609	23.609	
		CPU baseline:					
	CPI	283.989 ms J parallel baselir)A'				
	OI C	85.393 ms	ic.				

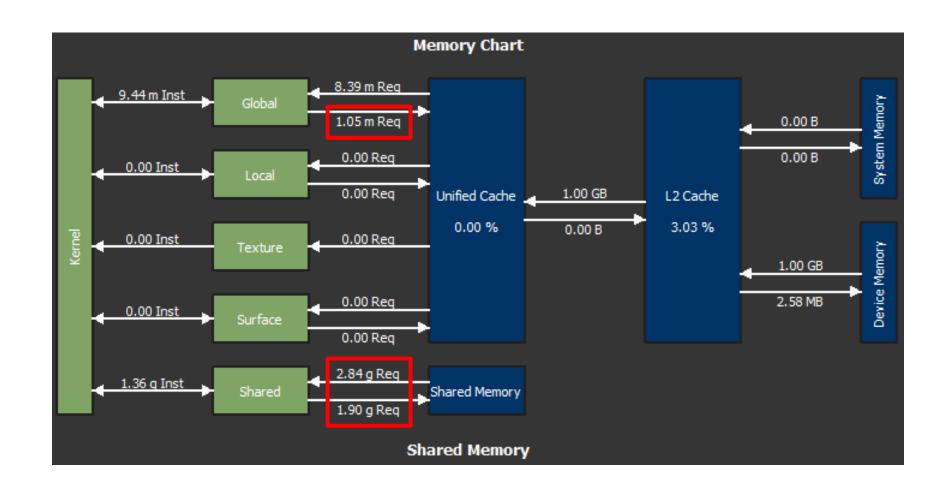
Parallel Reduction #2: Atomic Shared



	la stancation o	Democrate	9/ David	Dank Canffish						
C1 11 1	Instructions	Requests	% Peak	Bank Conflicts						
Shared Load Shared Store	404,404,197 8,388,608	2,842,279,716	98.40 65.89	1,020,076 1,028,529						
Shared Store Shared Atomic	8,388,608 943,634,854	1,903,184,212	65.89	1,028,529						
Total	1,356,427,659	4,745,463,928	164.29	2,048,605						
Total	1,530,427,039	4,743,403,926	104.29	2,046,003						
									First-Level Cad	ne
	Instructions	SM->TEX Requests	% Peak	Hit Rate TEX	->L2 Requests	% Peak	L2-> TEX Returns	% Peak	TEX->SM Returns	% Peak
Global Load Cached	8,388,608	8,388,608	0.29							
Global Load Uncached							33,554,432	1.16	1,905,957,978	65.99
Local Load Cached							55,55 4, 152		.,,,	33133
Local Load Uncached										
Surface Load									0	0
Texture Load										
Global Store					0	0				
Local Store										
	4 0 4 0 4 0 4	4.040.000	2		4.040.574	2				
Global Reduction	1,048,576	1,048,576	0.04	0	1,048,576	0.00				
Surrace Reduction Global Atomic	0	U	U	U	v	v				
									see above	see above
Global Atomic Cas Surface Atomic	0									
Surface Atomic Surface Atomic Cas	0								see above	see above
Loads	8,388,608	8,388,608	0.29	0			33,554,432	1.16	1,905,957,978	65.99
Stores	0,300,000	0,300,000	0.29	0	0	0	55,334,452	1.10	1,903,937,976	03.99
Total	9,437,184	9,437,184	0.33	0	1,048,576	0.00	33,554,432	1.16	1,905,957,978	65.99
Total	3,437,104	3,437,104	0.33	v	1,040,570	0.00	33,334,432	1.10		
									Second-Level Ca	che
	TEX->L2 Requests	% Peak L2->	TEX Returns	% Peak	Total Bytes	Total T	hroughput			
Global Load Cached					1,073,741,824		35,673,358,125.97			
Global Load Uncached			33,554,432	1.16						
Local Load Cached										
Local Load Uncached										
Surface Load							:			
Surface Load Texture Load						VS.	8.0 GB			
Surface Load Texture Load Global Store						VS.	8.0 GB			
Surface Load Texture Load						VS.	8.0 GB			
Surface Load Texture Load Global Store Local Store					0	VS.				
Surface Load Texture Load Global Store Local Store Control Control Global Reduction	1,048,576	5 0.00			0 2 33,554,432	VS.	8.0 GB			
Surface Load Texture Load Global Store Local Store					0	VS.				
Surface Load Texture Load Global Store Local Store Global Reduction Surface Reduction Global Atomic	1,048,576	5 0.00			0 2 33,554,432	VS.	1,114,792,441.44			
Surface Load Texture Load Global Store Local Store Global Reduction Surface Reduction Global Atomic Global Atomic Cas	1,048,576 0	0.00			33,554,432 U	VS.	2 1,114,792,441.44 v			
Surface Load Texture Load Global Store Local Store Global Reduction Juriace Reduction Global Atomic Global Atomic Cas Surface Atomic	1,048,576 0	0.00			33,554,432 U	VS.	2 1,114,792,441.44 v			
Surface Load Texture Load Global Store Local Store Global Reduction Surface Reduction Global Atomic Global Atomic Global Atomic Cas Surface Atomic Surface Atomic	1,048,576	0.00 0	0 0 - - - - 0	0 - - - 0 0	33,554,432 0 0	VS.	1,114,792,441.44 0 0			
Surface Load Texture Load Global Store Local Store Construction Global Reduction Surface Reduction Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads	1,048,576 0 0	0.00			0 33,554,432 0 0 0	VS.	1,114,792,441.44 0 0 0 35,673,358,125.97			
Surface Load Texture Load Global Store Local Store Global Reduction surrace neduction Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores	1,048,576 0 0	0.00	0 0 - - - 0 0 33,554,432	0 - - - 0 0	0 2 33,554,432 0 0 0 1,073,741,824 0	VS.	0 0 0 0 35,673,358,125.97			
Surface Load Texture Load Global Store Local Store Construction Global Reduction Surface Reduction Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads	1,048,576 0 0	0.00	0 0 - - - - 0	0 - - - 0 0	0 33,554,432 0 0 0	VS.	1,114,792,441.44 0 0 0 35,673,358,125.97			
Surface Load Texture Load Global Store Local Store Comment of the Store Global Reduction Surface Reduction Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores	1,048,576 0 0	0.00	0 0 - - - 0 0 33,554,432	0 - - - 0 0	0 2 33,554,432 0 0 0 1,073,741,824 0	VS.	0 0 0 0 35,673,358,125.97		Device Memor	y
Surface Load Texture Load Global Store Local Store Comment of the Store Global Reduction Surface Reduction Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores	1,048,576 0 0 0 - 0 1,048,576	0.00 0 0 0 - 0 0.00	0 0 - - 0 0 33,554,432 - 33,554,432	0 - - - 0 0 1.16 - 1.16	33,554,432 0 0 1,073,741,824 0 1,107,296,256	VS.	0 0 0 0 35,673,358,125.97		Device Memor	y
Surface Load Texture Load Global Store Local Store Global Reduction surrace neduction Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores	1,048,576 0 0 0 - 0 1,048,576	0.00	0 0 - - 0 0 33,554,432 - 33,554,432	0 - - - 0 0 1.16 - 1.16	33,554,432 0 0 0 1,073,741,824 0 1,107,296,256	VS.	0 0 0 0 35,673,358,125.97		Device Memor	y
Surface Load Texture Load Global Store Local Store Global Reduction surrace neauction Global Atomic Global Atomic Cas Surface Atomic Surface Atomic Cas Loads Stores Total	1,048,576 0 0 0 - 0 1,048,576	0.00 0 0 0 - 0 0.00	0 0 - - 0 0 33,554,432 - 33,554,432	0 - - - 0 0 1.16 - 1.16 Throughpi	33,554,432 0 0 1,073,741,824 0 1,107,296,256	VS.	0 0 0 0 35,673,358,125.97		Device Memor	У

Parallel Reduction #2: Atomic Shared



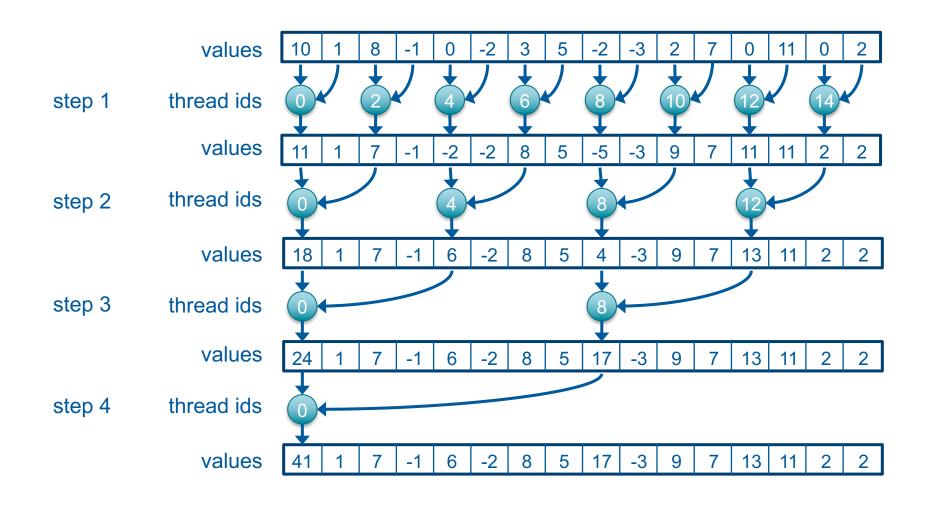




```
global__ void reduce_shared(const float* input, float* result, int elements) {
 extern __shared__ float data[];
 int id = threadIdx.x + blockIdx.x*blockDim.x;
  data[threadIdx.x] = input[id];
  __syncthreads();
 for (int s = 1; s < blockDim.x; s *= 2) {</pre>
      if (threadIdx.x % (2*s) == 0)
          data[threadIdx.x] += data[threadIdx.x + s];
       _syncthreads();
 if (threadIdx.x == 0)
      atomicAdd(result, data[0]);
```

Parallel Reduction #3: Shared



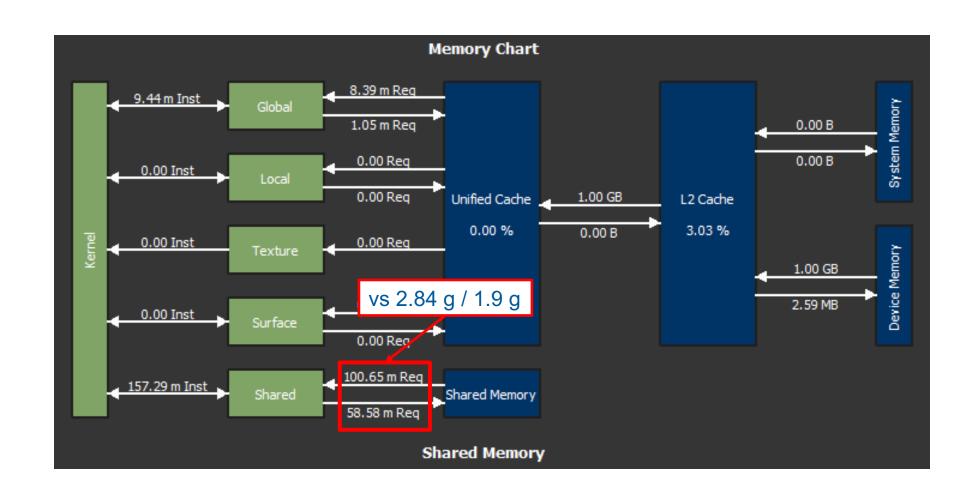




	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000

Parallel Reduction #3: Shared



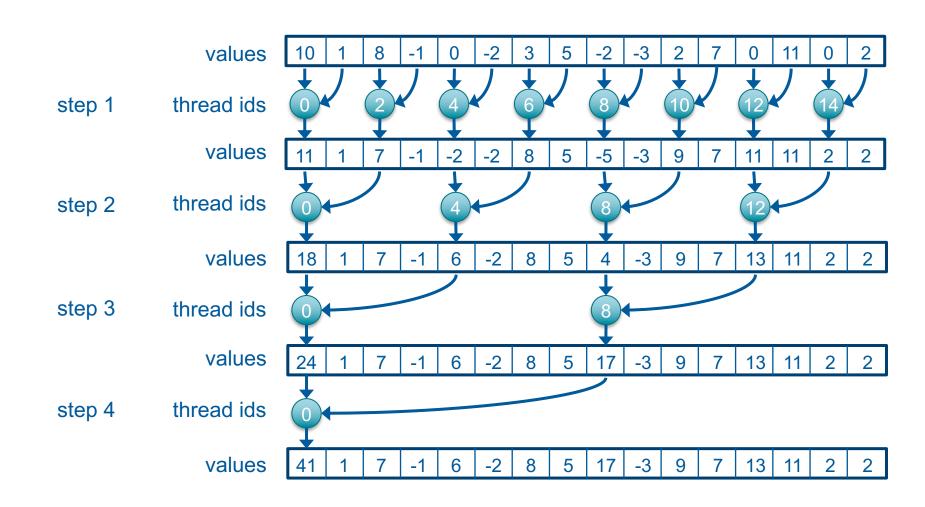


Parallel Reduction #3: Shared





Divergence



Parallel Reduction #4: Shared Anti Divergence

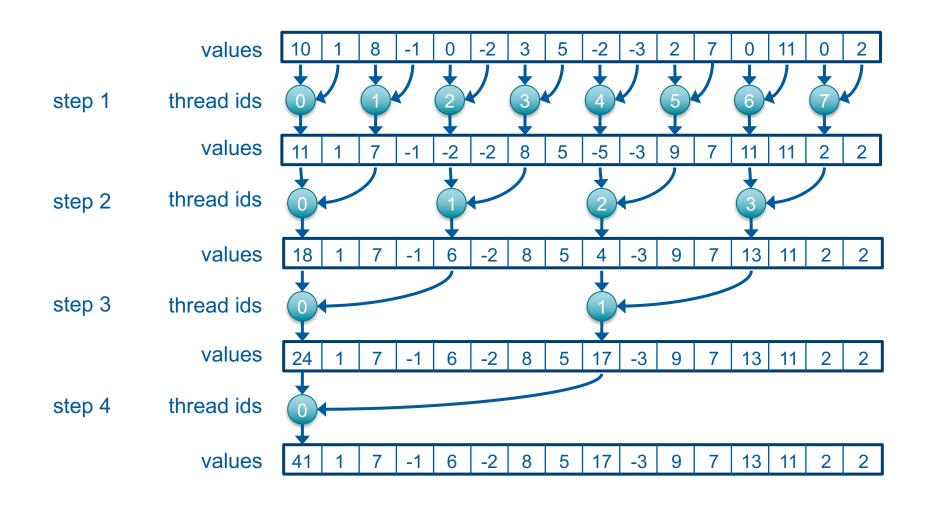


```
for (int s = 1; s < blockDim.x; s *= 2) {
   if (threadIdx.x % (2*s) == 0)
      data[threadIdx.x] += data[threadIdx.x + s];
   __syncthreads();
}</pre>
```

```
for (int s = 1; s < blockDim.x; s *= 2) {
   int lid = 2*s*threadIdx.x;
   if (lid < blockDim.x)
       data[lid] += data[lid + s];
   __syncthreads();
}</pre>
```

Parallel Reduction #4: Shared Anti Divergence







	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000
anti divergence	3.849 ms	250.859 GB/s	1.630	165.070	1.630

Parallel Reduction #4: Shared Anti Divergence



Problem

Bank conflicts

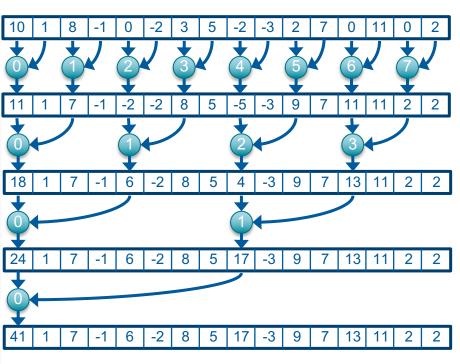
	Instructions	Requests	% Peak	Bank Conflicts
Shared Load	26,214,400	101,061,582	24.55	73,419,745
Shared Store	20,971,520	58,645,681	14.25	37,674,161
Shared Atomic	0	-	-	-
Total	47,185,920	159,707,263	38.80	111,093,906

If multiple addresses of a memory request map to the same memory bank, the accesses are serialized! Bank conflicts can occur between any threads in the warp!

Parallel Reduction #5: Shared Anti Conflicts

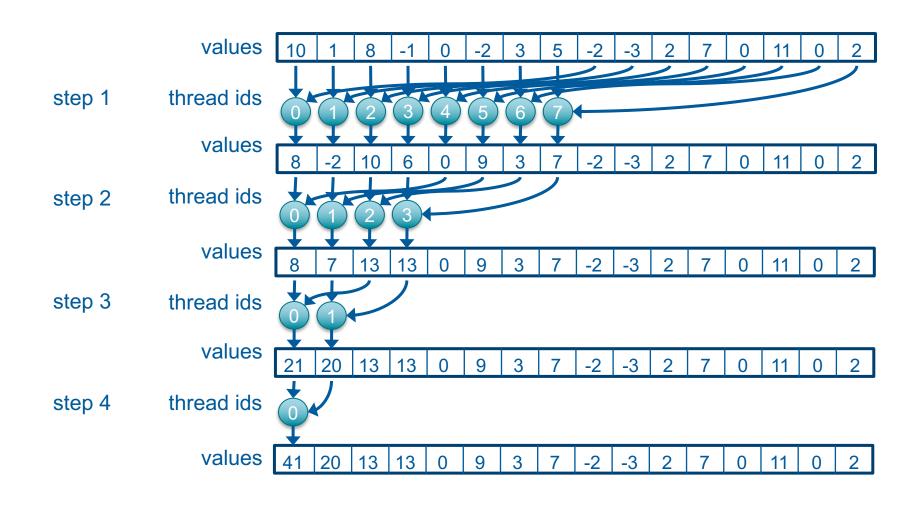


```
for (int s = 1; s < blockDim.x; s *= 2) {</pre>
    int lid = 2*s*threadIdx.x;
   if (lid < blockDim.x)</pre>
        data[lid] += data[lid + s];
    __syncthreads();
for (int s = blockDim.x / 2; s > 0; s /= 2) {
   if (threadIdx.x < s)</pre>
        data[threadIdx.x] += data[threadIdx.x + s];
    __syncthreads();
```



Parallel Reduction #5: Shared Anti Conflicts







	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000
anti divergence	3.849 ms	250.859 GB/s	1.630	165.070	1.630
anti conflicts	2.903 ms	333.536 GB/s	1.326	218.862	2.162

Parallel Reduction #5: Shared Anti Conflicts



	Instructions	Requests	% Peak	Bank Conflicts
Shared Load	26,214,400	26,924,195	8.71	63,568
Shared Store	20,971,520	21,875,836	7.08	904,316
Shared Atomic	0	-	-	-
Total	47,185,920	48,800,031	15.79	967,884
				vs 111,093,906

Parallel Reduction #5: Shared Anti Conflicts



- Problem
 - Half of threads are idle

```
for (int s = blockDim.x / 2; s > 0; s /= 2) {
   if (threadIdx.x < s)
      data[threadIdx.x] += data[threadIdx.x + s];
   __syncthreads();
}</pre>
```



```
global__ void reduce_shared_dual_load(const float* input, float* result, int elements) {
 extern __shared float data[];
  int id = threadIdx.x + blockIdx.x*blockDim.x * 2;
  data[threadIdx.x] = input[id] + input[id + blockDim.x];
  syncthreads();
 for (int s = blockDim.x / 2; s > 0; s /= 2) {
      if (threadIdx.x < s)</pre>
          data[threadIdx.x] += data[threadIdx.x + s];
      _syncthreads();
 if (threadIdx.x == 0)
      atomicAdd(result, data[0]);
```

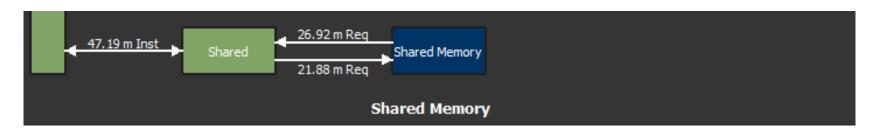


	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000
anti divergence	3.849 ms	250.859 GB/s	1.630	165.070	1.630
anti conflicts	2.903 ms	333.536 GB/s	1.326	218.862	2.162
dual load	1.66 ms	612.692 GB/s	1.749	382.744	3.780

Parallel Reduction #6: Shared Dual Load



shared_anti_conflicts<<<1048576, 256>>>(...)



shared_dual_load<<<524288, 256>>>(...)

```
23.59 m Inst
Shared
Shared Memory
11.18 m Req
Shared Memory
```

Parallel Reduction #6: Shared Dual Load



- 612 GB/s vs 651 GB/s peak
- What is the bottleneck?
 - Probably instruction overhead
 - Operations that are not loads, stores, or computations
 - Address arithmetic, loops, conditionals
- Idea: unroll loops

Parallel Reduction #7: Shared Unrolled Last



```
extern __shared__ volatile float data[];
for (int s = blockDim.x / 2; s > 32; s /= 2) {
   if (threadIdx.x < s)</pre>
        data[threadIdx.x] += data[threadIdx.x + s];
    syncthreads();
if (threadIdx.x < 32) {</pre>
                                                                      CC \geq 7.0
    data[threadIdx.x] += data[threadIdx.x + 32];
                                                          Watch out:
    data[threadIdx.x] += data[threadIdx.x + 16];
    data[threadIdx.x] += data[threadIdx.x + 8];
                                                     With ITS this requires
    data[threadIdx.x] += data[threadIdx.x + 4];
                                                      syncwarp() after
    data[threadIdx.x] += data[threadIdx.x + 2];
                                                          each step!
    data[threadIdx.x] += data[threadIdx.x + 1];
  (threadIdx.x == 0)
    atomicAdd(result, data[0]);
```

Parallel Reduction #7: Shared Unrolled Last



```
extern __shared__ float data[];
float x = 0.0f;
for (int s = blockDim.x / 2; s > 32; s /= 2) {
    if (threadIdx.x < s)</pre>
        data[threadIdx.x] += data[threadIdx.x + s];
     _syncthreads();
if (threadIdx.x < 32) {</pre>
    x += data[threadIdx.x + 32]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 16]; syncwarp(); data[threadIdx.x] = x; syncwarp();
    x += data[threadIdx.x + 8]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 4]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 2]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 1]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
if (threadIdx.x == 0)
    atomicAdd(result, data[0]);
```



	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000
anti divergence	3.849 ms	250.859 GB/s	1.630	165.070	1.630
anti conflicts	2.903 ms	333.536 GB/s	1.326	218.862	2.162
dual load	1.66 ms	612.692 GB/s	1.749	382.744	3.780
unrolled last	1.65 ms	646.565 GB/s	1.006	385.063	3.803

Parallel Reduction #8: Shared Unrolled All (1/2)



```
template<int BlockSize>
 global void reduce shared unrolled_all(const float* input, float* result, int elements) {
   extern __shared__ float data[];
    int id = threadIdx.x + blockIdx.x*blockDim.x * 2;
    data[threadIdx.x] = input[id] + input[id + blockDim.x];
     _syncthreads();
    if (BlockSize >= 512) {
        if (threadIdx.x < 256)</pre>
            data[threadIdx.x] += data[threadIdx.x + 256];
        __syncthreads();
       (BlockSize >= 256) {
        if (threadIdx.x < 128)</pre>
            data[threadIdx.x] += data[threadIdx.x + 128];
         _syncthreads();
```

Parallel Reduction #8: Shared Unrolled All (2/2)



```
if (BlockSize >= 128) {
   if (threadIdx.x < 64)</pre>
        data[threadIdx.x] += data[threadIdx.x + 64];
    __syncthreads();
if (threadIdx.x < 32) {</pre>
    x += data[threadIdx.x + 32]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 16]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 8]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 4]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 2]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 1]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
if (threadIdx.x == 0)
    atomicAdd(result, data[0]);
```



	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000
anti divergence	3.849 ms	250.859 GB/s	1.630	165.070	1.630
anti conflicts	2.903 ms	333.536 GB/s	1.326	218.862	2.162
dual load	1.66 ms	612.692 GB/s	1.749	382.744	3.780
unrolled last	1.65 ms	646.565 GB/s	1.006	385.063	3.803
unrolled all	1.65 ms	646.761 GB/s	1.000	385.063	3.803

Parallel Reduction #9: Shared Shuffle



```
if (threadIdx.x < 32) {
    x += data[threadIdx.x + 32]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 16]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 8]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 4]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 2]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
    x += data[threadIdx.x + 1]; __syncwarp(); data[threadIdx.x] = x; __syncwarp();
}</pre>
```

```
if (threadIdx.x < 32) {
    x = data[threadIdx.x] + data[threadIdx.x + 32];
    x += __shfl_sync(@xFFFFFFFFF, x, threadIdx.x + 16);
    x += __shfl_sync(@x00000FFFF, x, threadIdx.x + 8);
    x += __shfl_sync(@x0000000FF, x, threadIdx.x + 4);
    x += __shfl_sync(@x00000000FF, x, threadIdx.x + 2);
    x += __shfl_sync(@x0000000000FF, x, threadIdx.x + 2);
}</pre>
```



	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000
anti divergence	3.849 ms	250.859 GB/s	1.630	165.070	1.630
anti conflicts	2.903 ms	333.536 GB/s	1.326	218.862	2.162
dual load	1.66 ms	612.692 GB/s	1.749	382.744	3.780
unrolled last	1.65 ms	646.565 GB/s	1.006	385.063	3.803
unrolled all	1.65 ms	646.761 GB/s	1.000	385.063	3.803
shuffle	1.65 ms	647.059 GB/s	1.000	385.063	3.803

Reduction Complexity



CC ≥ 3.0

It is possible to directly exchange registers within a warp

__shfl_sync(unsigned mask, float val, int src, int width=WarpSize)

	Instructions	Requests	% Peak	Bank Conflicts
Shared Load	7,340,032 (-41.67%)	7,470,212 (-41.86%)	3.60 (-41.56%)	16,636 (-52.02%)
Shared Store	7,340,032 (-26.32%)	7,929,460 (-26.15%)	3.82 (-25.77%)	589,428 (-24.01%)
Shared Atomic	0 (+0.00%)	-	-	-
Total	14,680,064 (-34.88%)	15,399,672 (-34.71%)	7.41 (-34.37%)	606,064 (-25.21%)

Parallel Reduction #10: Shared Multi Seq



```
extern __shared__ float data[];
int id = threadIdx.x + blockIdx.x*blockDim.x * 2;
data[threadIdx.x] = input[id] + input[id + threadIdx.x];
__syncthreads();
```

```
extern __shared__ float data[];
int id = threadIdx.x + blockIdx.x*blockDim.x;
int grid_size = BlockSize*gridDim.x;

float in = 0.0f;
for (; id < elements; id += grid_size)
    in += input[id];

data[threadIdx.x] = in;
__syncthreads();</pre>
```

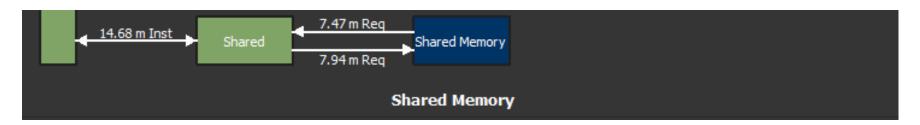


	Time (2 ²⁸ floats)	Bandwidth	Step speedup	Cumulative speedup	Speedup to reduction
atomic global	635.355 ms	1.524 GB/s	1.000	1.000	
atomic shared	26.911 ms	35.764 GB/s	23.609	23.609	
reduction	6.275 ms	153.726 GB/s	4.288	101.252	1.000
anti divergence	3.849 ms	250.859 GB/s	1.630	165.070	1.630
anti conflicts	2.903 ms	333.536 GB/s	1.326	218.862	2.162
dual load	1.66 ms	612.692 GB/s	1.749	382.744	3.780
unrolled last	1.65 ms	646.565 GB/s	1.006	385.063	3.803
unrolled all	1.65 ms	646.761 GB/s	1.000	385.063	3.803
shuffle	1.65 ms	647.059 GB/s	1.000	385.063	3.803
multi seq	1.69 ms	641.734 GB/s	0.976	375.950	3.713

Parallel Reduction #10: Shared Multi Seq



shared_shuffle<<<524288, 256>>>(...)



shared_multi_seq<<<37450, 256>>>(...)

```
524.48 k Req
1.05 m Inst
                 Shared
                                          Shared Memory
                             557.01 k Req
                                       Shared Memory
```

Reduction Performance: Different Architectures



	2 ²⁸ GTX 780	2 ²⁸ GTX 970	2 ²⁸ GTX 1080 Ti	2 ²⁸ Titan V	2 ²⁸ RTX 2080 Ti
atomic global	893.79 ms	644.64 ms	451.37 ms	635.35 ms	437.92 ms
atomic shared	558.49 ms	3433.07 ms	649.13 ms	26.91 ms	47.46 ms
reduction	56.72 ms	16.78 ms	17.76 ms	6.28 ms	7.41 ms
anti divergence	42.16 ms	38.50 ms	10.62 ms	3.84 ms	6.02 ms
anti conflicts	31.66 ms	30.45 ms	8.80 ms	2.90 ms	4.83 ms
dual load	17.68 ms	14.35 ms	4.58 ms	1.66 ms	2.53 ms
unrolled last	11.93 ms	10.18 ms	3.08 ms	1.65 ms	1.87 ms
unrolled all	10.37 ms	8.43 ms	3.06 ms	1.65 ms	1.84 ms
shuffle	9.89 ms	8.04 ms	3.03 ms	1.65 ms	1.83 ms
multi seq	7.53 ms	6.95 ms	2.99 ms	1.69 ms	1.87 ms