

[Total Recall – Columbia Pictures]

CUDA Memory Spaces

CUDA GPU with Caches



Memory spaces and performance

- Global Memory
 - cached (L2 and L1)
 - read and write
 - dynamic allocation in device code possible
- Constant Memory
 - cached
 - read-only
- Texture Memory
 - cached
 - interpolation
 - ‘read-only’
- Shared Memory
 - communication between threads
- Register
 - private for each thread
 - (exchangeable)

Global Memory

- Main data storage
 - Use for input and output data
 - Linear arrays
- Access pattern important
- Cache behavior important
- Relatively slow
 - Bandwidth: $\sim 180\text{-}280\text{ GB/s}$
 - Non-cached coalesced access: 300 cycles/ns
 - L2 cached access: 160 cycles/ns
 - L1 caches access: 20 cycles/ns

Cache usage

- Not intended for the same use as CPU caches
 - Smaller size (especially per thread)
 - Not aimed for temporal reuse
 - Smooth-out access patterns
 - help with spilled registers (L1 on Kepler)
- Don't try cache blocking like on CPU
 - 100s of threads accessing L1
 - 1000s of threads accessing L2
 - use shared memory instead
- Optimizations should not target the cache architecture

Global memory operations

- Cacheline Size
 - L1: 128 Byte
 - L2: 32 Byte
- Stores
 - Write-through for L1
 - Write-back for L2
- Memory operations are issued per warp
 - Threads provide addresses
 - Combined to lines/segments needed
 - Requested and served
- Try to get coalescing per warp
 - Align starting address
 - Access within a contiguous region

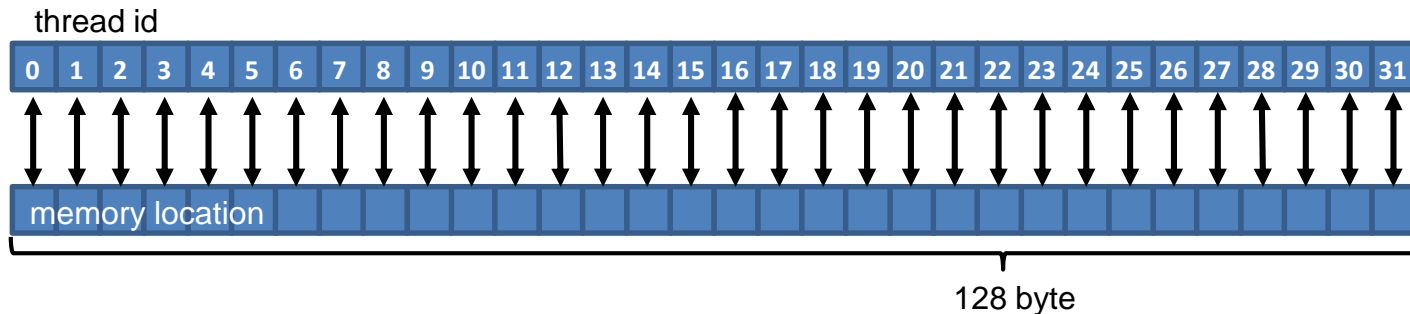
Granularity Example 1/4

```

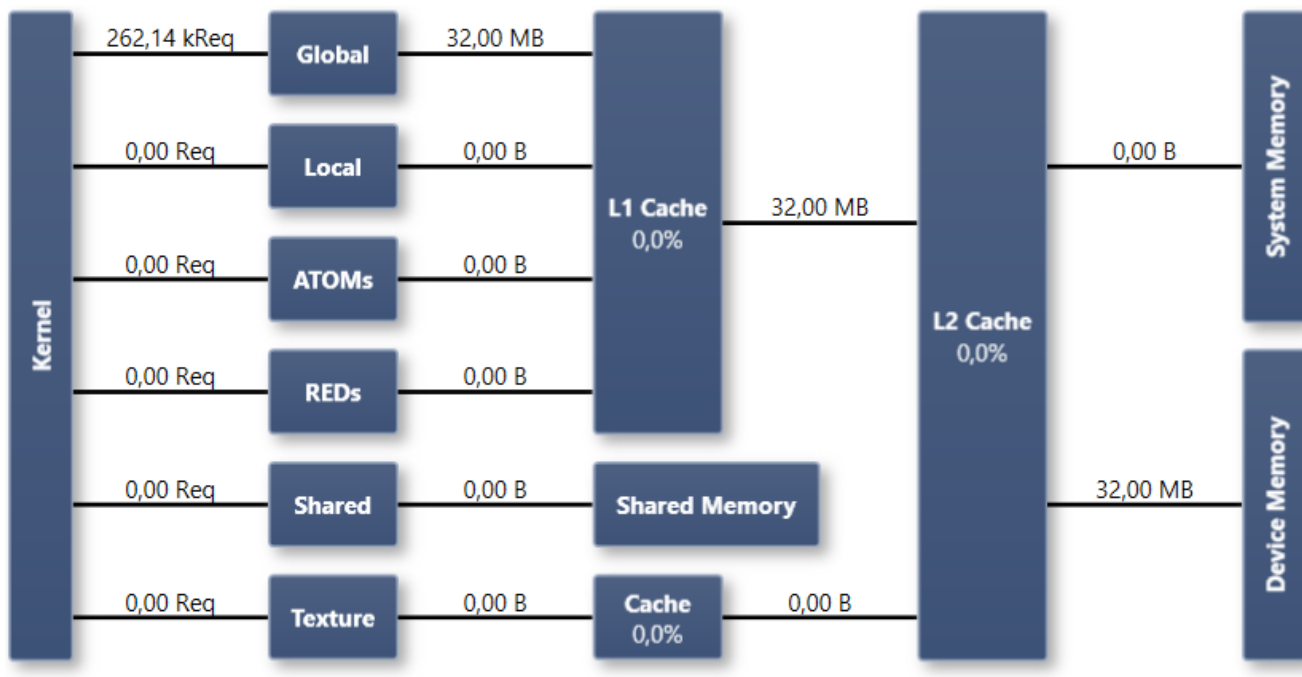
7 | __global__ void testCoalesced(int *in, int* out, int elements)
8 | {
9 |   int block_offset = blockIdx.x*blockDim.x;
10 |  int warp_offset = 32*(threadIdx.x/32);
11 |  int laneid = threadIdx.x%32;
12 |  int id = (block_offset + warp_offset + laneid)%elements;
13 |
14 |  out[id] = in[id];
15 | }

```

testCoalesced done in 0.240736ms -> 69.6914GB/s



Granularity Example 1/4



Coalesced

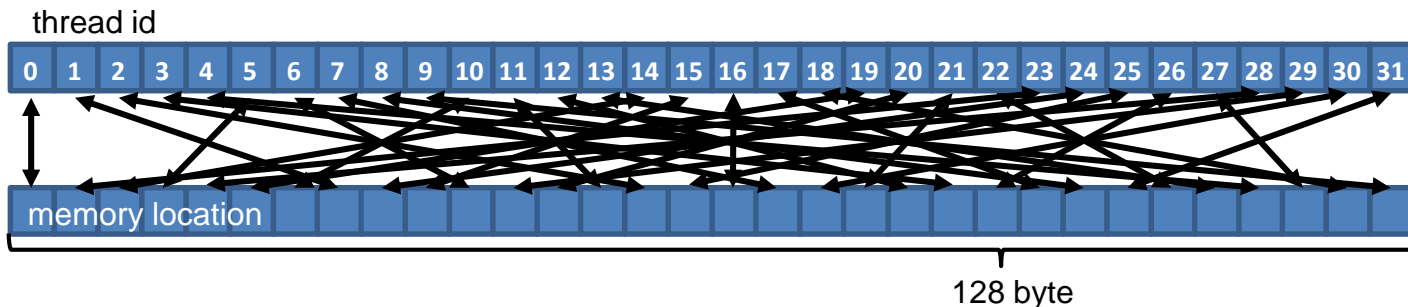
Granularity Example 2/4

```

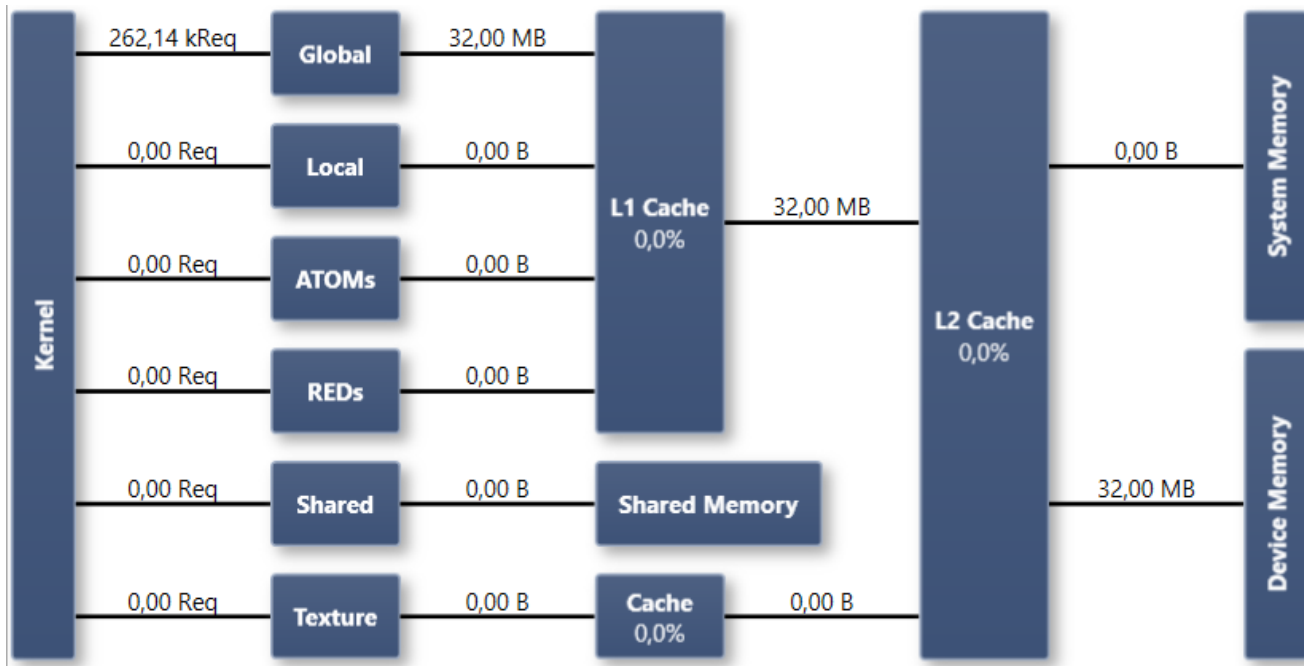
16 __global__ void testMixed(int *in, int* out, int elements)
17 {
18   int block_offset = blockIdx.x*blockDim.x;
19   int warp_offset = 32*(threadIdx.x/32);
20   int elementid = (threadIdx.x*7)%32;
21   int id = (block_offset + warp_offset + elementid)%elements;
22
23   out[id] = in[id];
24 }

```

testCoalesced done in 0.240736ms -> 69.6914GB/s
 testMixed done in 0.24736ms -> 67.8251GB/s



Granularity Example 2/4



Mixed

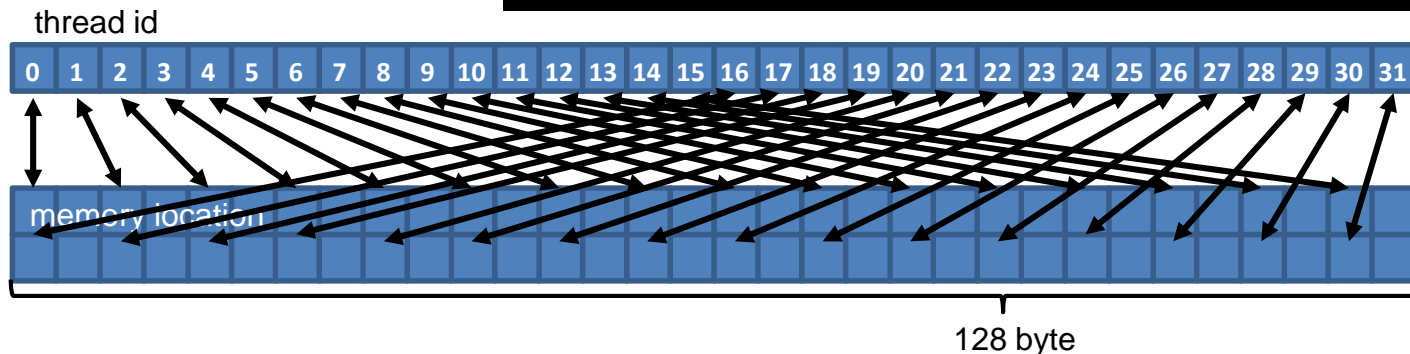
Granularity Example 3/4

```

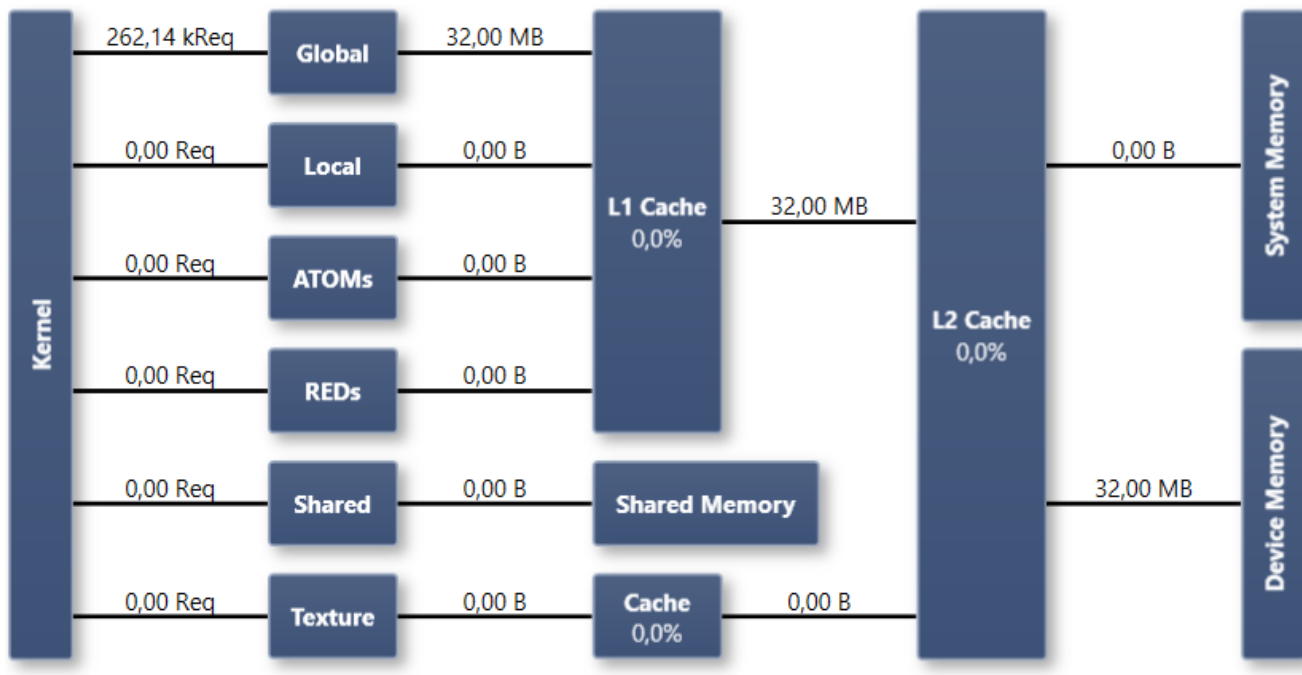
23  template<int offset>
24  □ __global__ void testOffset(int *in, int* out, int elements)
25  {
26      int block_offset = blockIdx.x*blockDim.x;
27      int warp_offset = 32*(threadIdx.x/32);
28      int laneid = threadIdx.x%32;
29      int id = ((block_offset + warp_offset + laneid)*offset)%elements;
30
31      out[id] = in[id];|
32  }

```

testCoalesced -> 69.69GB/s
 testOffset<2> -> 40.80GB/s testOffset<4> -> 20.27GB/s
 testOffset<8> -> 10.06GB/s testOffset<32> -> 6.54GB/s

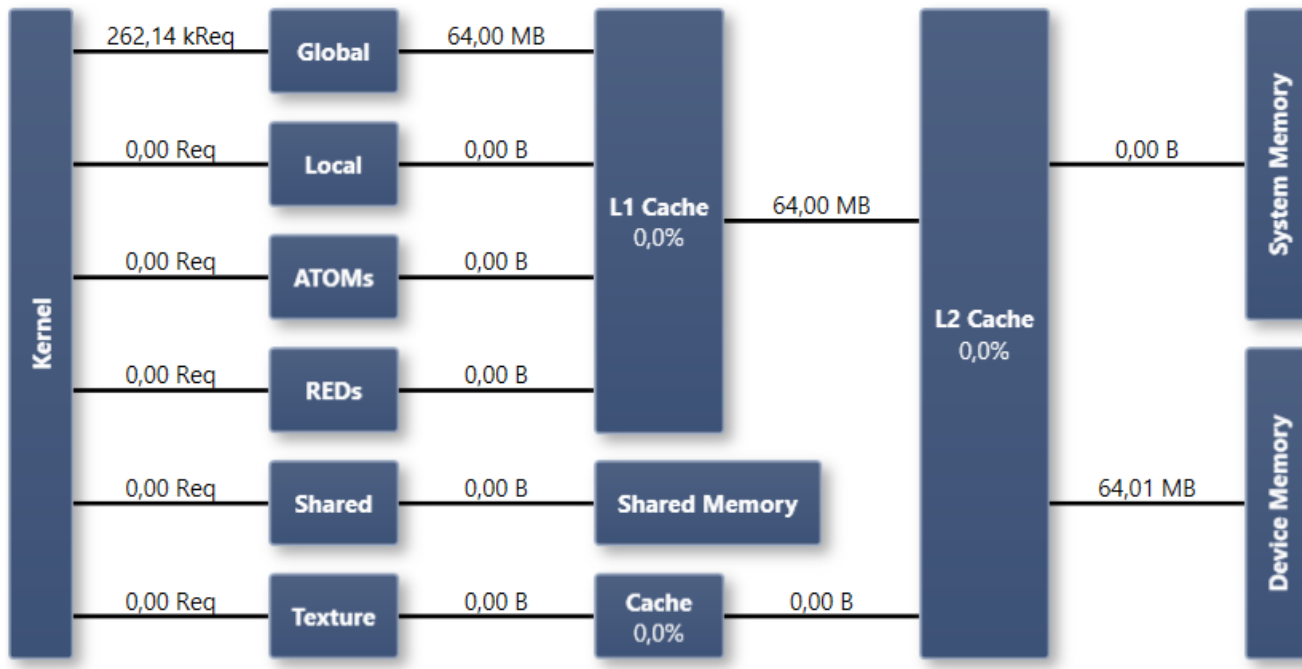


Granularity Example 3/4



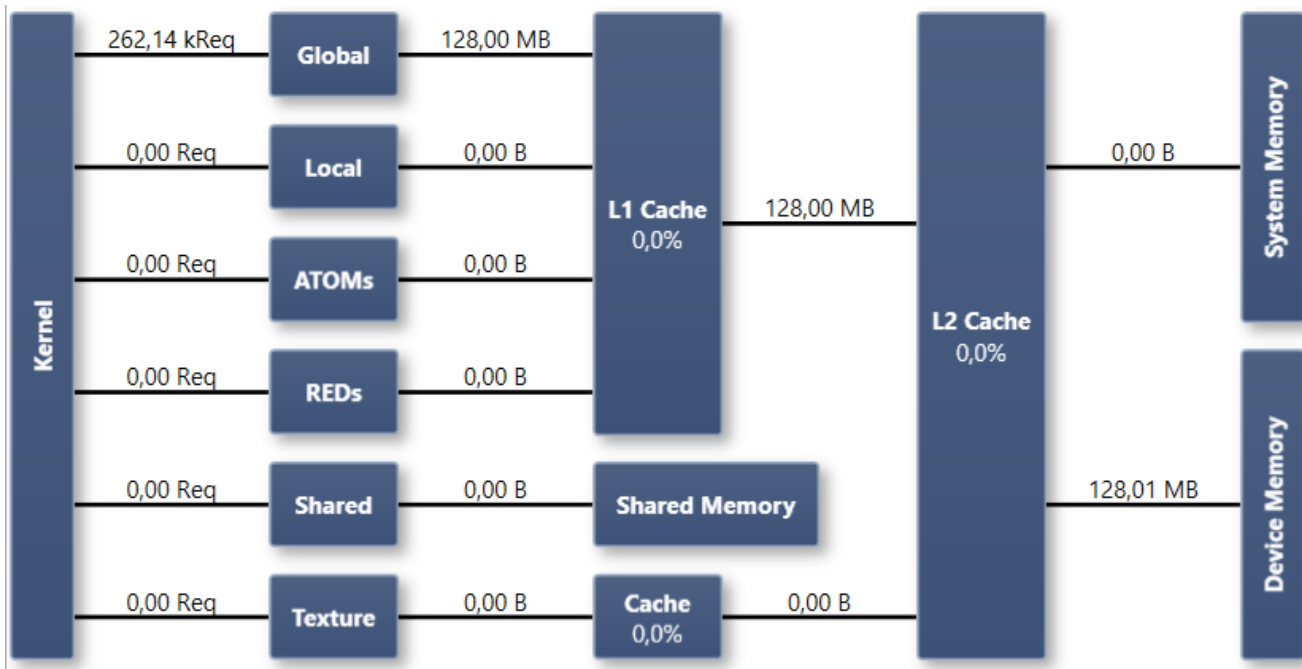
Coalesced 69.69GB/s

Granularity Example 3/4



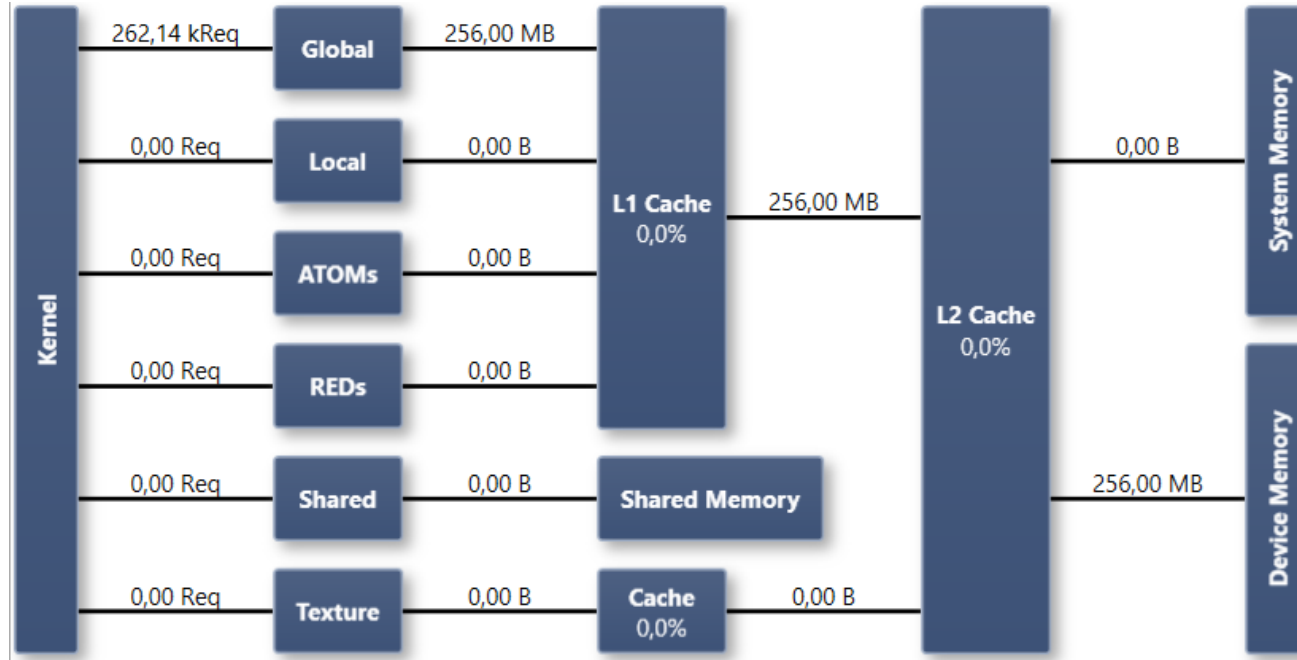
Offset 2 40.80GB/s (0.58)

Granularity Example 3/4



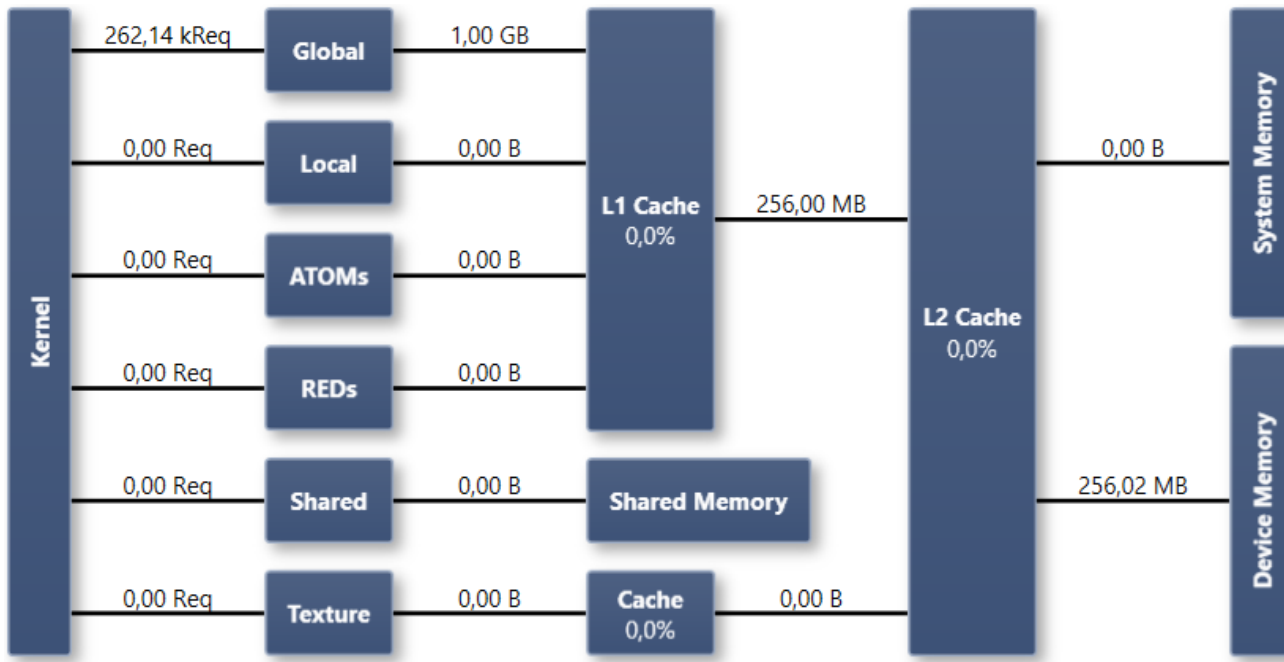
Offset 4 20.27GB/s (0.29)

Granularity Example 3/4



Offset 8 10.06GB/s (0.14)

Granularity Example 3/4



Offset 32 6.54GB/s (0.09)

Granularity Example 4/4

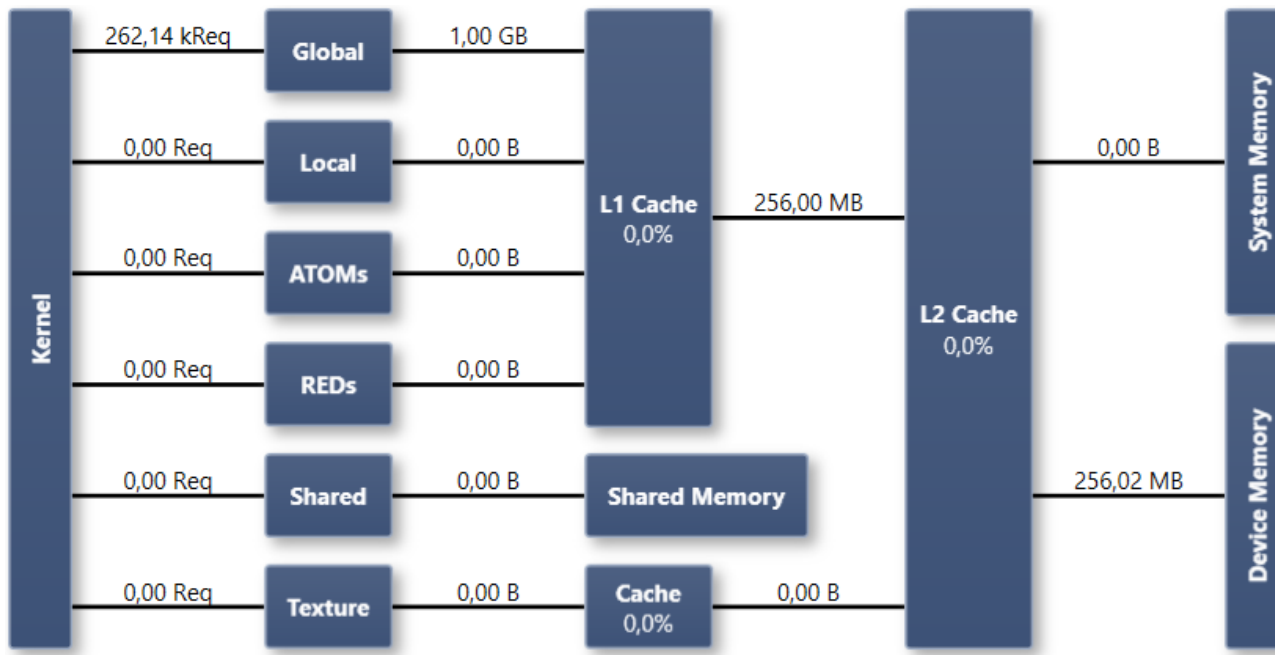
```

36 __global__ void testScattered(int *in, int* out, int elements)
37 {
38     int block_offset = blockIdx.x*blockDim.x;
39     int warp_offset = 32*(threadIdx.x/32);
40     int elementid = threadIdx.x%32;
41     int id = ((block_offset + warp_offset + elementid)*121)%elements;
42
43     out[id] = in[id];
44 }

```

testCoalesced done in 0.240736ms -> 69.6914GB/s
 testOffset<32> done in 2.5617ms -> 6.54351GB/s
 testScattered done in 4.70234ms -> 3.56785GB/s

Granularity Example 4/4



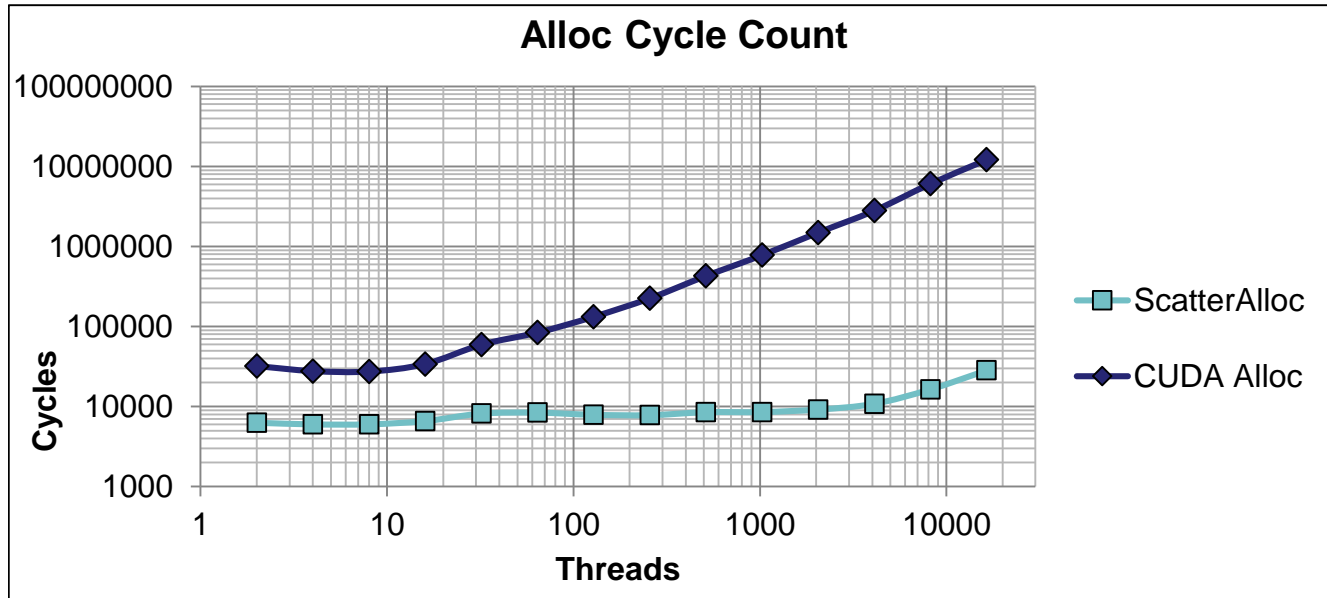
Scattered 3.57GB/s (0.05)

Granularity Example Summary

- testCoalesced done in 0.2407ms -> 69.6914GB/s
 - testMixed done in 0.2478ms -> 67.6937GB/s
 - testOffset<2> done in 0.4111ms -> 40.8073GB/s
 - testOffset<4> done in 0.8276ms -> 20.2717GB/s
 - testOffset<8> done in 1.6669ms -> 10.0647GB/s
 - testOffset<32> done in 2.5617ms -> 6.54351GB/s
 - testScattered done in 4.7023ms -> 3.56785GB/s
-
- Access pattern within 128 Byte segment does not matter
 - Offset between data → more requests need to be handled
 - Peak performance not met due to computation overhead
 - More scattered data access slower with GDDR RAM

Dynamic Memory Allocation

- malloc / free, new / delete are possible for cc ≥ 2.0
- Very slow 20 000 – 10 000 000 cycles



Constant memory

- Read-only
- Ideal for coefficients and other data that is read uniformly by warps
- Data is stored in global memory, read through cache
 - `__constant__` qualifier
 - limited to 64KB
- Fermi added uniform access:
 - Load Uniform (LDU) command for const pointer
 - Goes through same cache
 - Compiler determines that all threads in a block will dereference the same address
 - no size limit

Constant memory cont'd

- Cache throughput
 - 4 bytes per warp per clock
 - All threads in warp read the same address
 - Otherwise serialized

```

__constant__ float myarray[128];
__global__ void kernel()
{
    ...
    float x = myarray[23];           //uniform
    float y = myarray[blockIdx.x + 2]; //uniform
    float z = myarray[threadIdx.x];  //non-uniform
}

```

Constant Memory Example

```

__constant__ float c_array[128];
__device__ float d_array[128];
__device__ const float dc_array[128];
__global__ void kernel()
{
    float a = c_array[0];           // ld.const    25ns
    float b = c_array[blockIdx.x]; // ld.const    25ns
    float c = c_array[threadIdx.x]; // ld.const   920ns
    float d = d_array[blockIdx.x];  // ld.global   21ns
    float e = d_array[threadIdx.x]; // ld.global   42ns
    float f = dc_array[blockIdx.x]; // ldu         21ns
    float g = dc_array[threadIdx.x]; // ld.global   42ns
}

```

Constant Memory Example 2

```

__constant__ float c_array[1024];
__device__ float d_array[1024];
__device__ const float dc_array[1024];
__global__ void kernel()
{
    for(int i = 0; i < 800; i+= 8)
    {
        float a = c_array[i];           // ld.const    25ns
        float b = c_array[blockIdx.x+i]; // ld.const    25ns
        float c = c_array[threadIdx.x+i]; // ld.const  1195ns
        float d = d_array[blockIdx.x+i]; // ld.global   84ns
        float e = d_array[threadIdx.x+i]; // ld.global  127ns
        float f = dc_array[blockIdx.x+i]; // ldu         84ns
        float g = dc_array[threadIdx.x+i]; // ld.global  127ns
    }
}

```

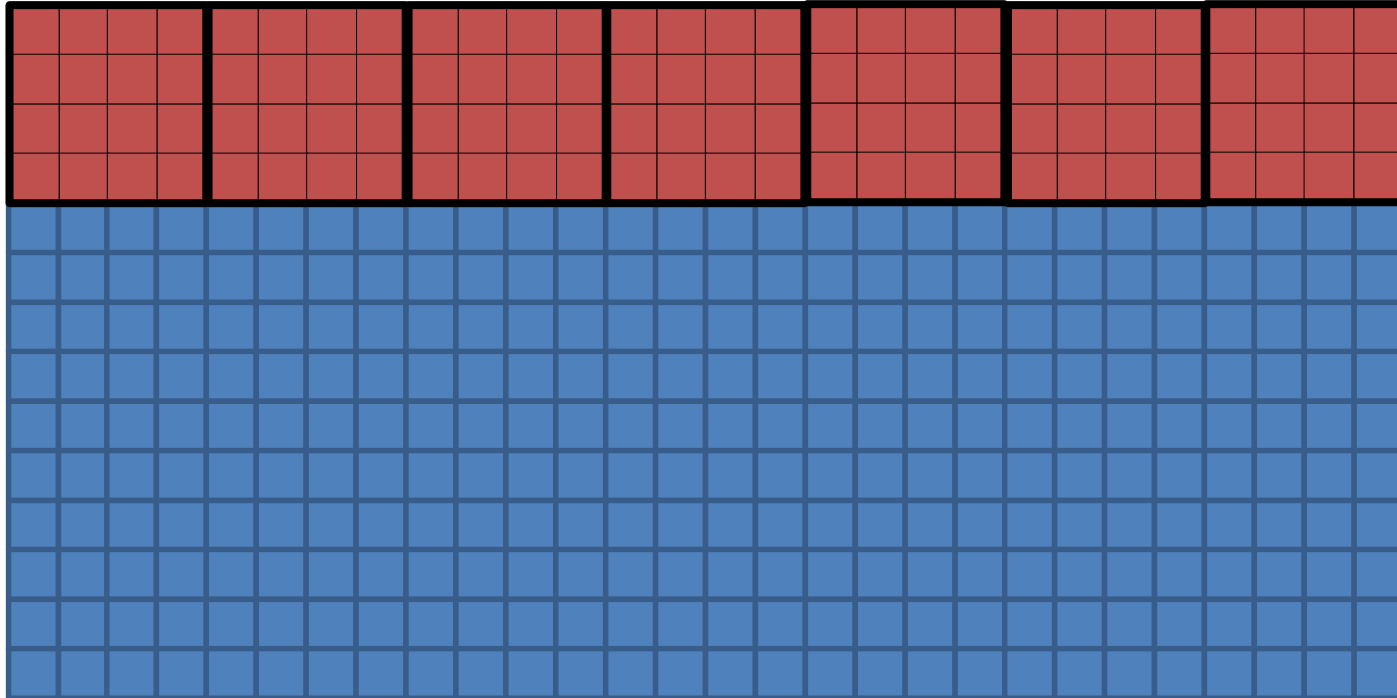

Constant memory

- Only useable if all threads within a warp read the same value
- Constant can be faster than global
- Uses different cache than global
- Compiler can automatically put things to constant

Texture memory

- Separate cache
- Textures are laid out in memory to optimize data locality for the respective dimensionality (cudaArray)
 - e.g. accesses to 2D region should hit the cache
- Surfaces allow writing to cudaArrays
- Concurrent writing and reading undefined result
- Can also bind global memory to textures

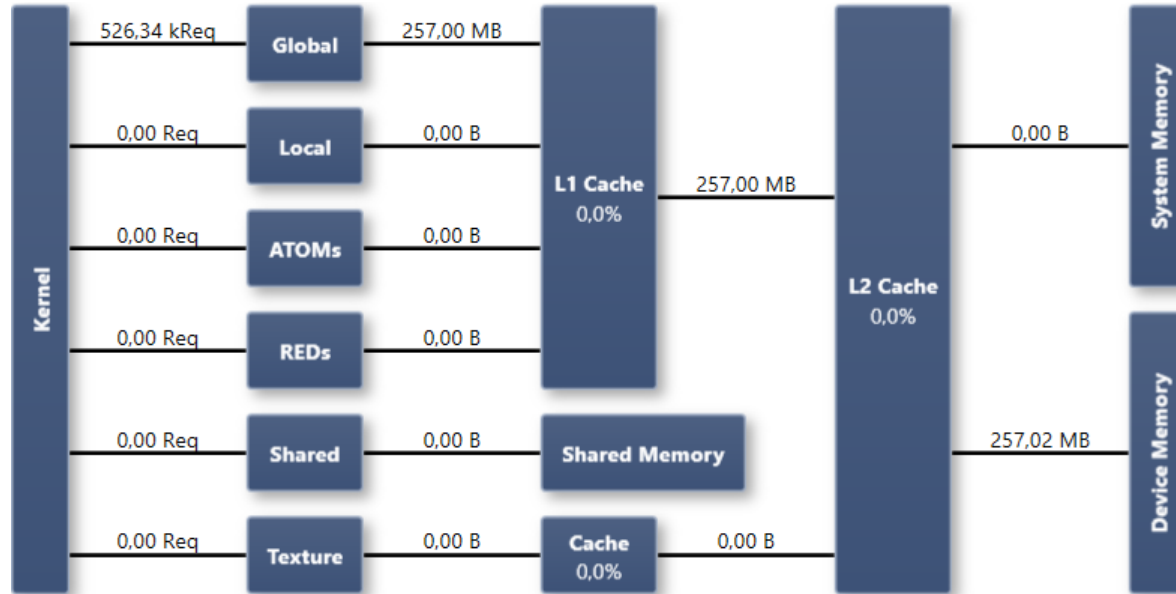
RowAfterRow



RowAfterRow: Texture vs Global

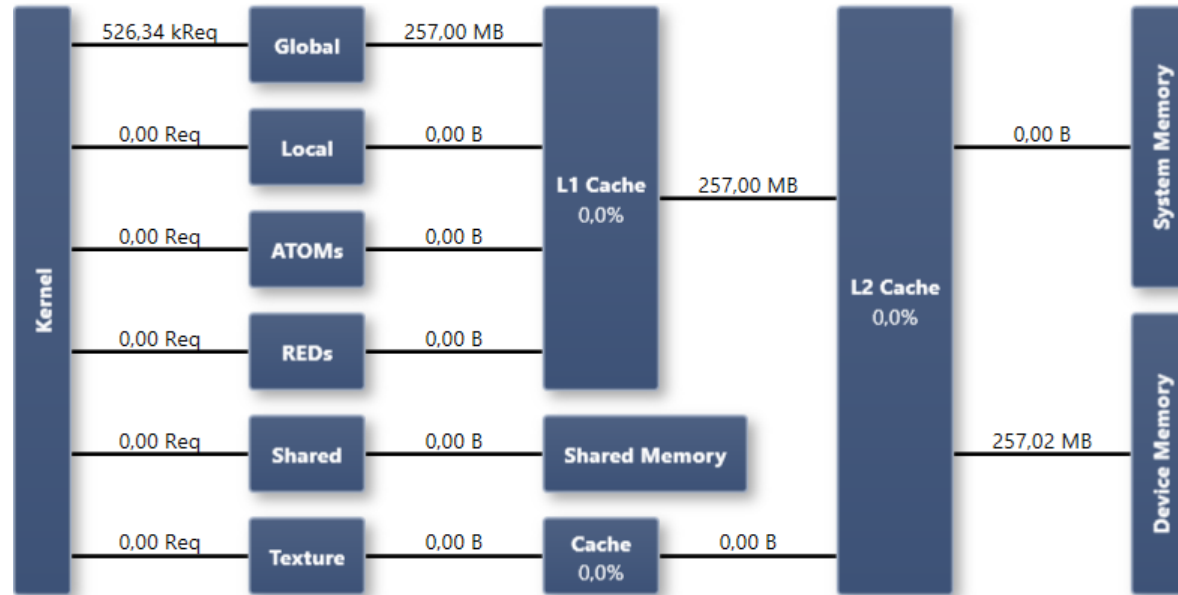
```
7  __global__ void deviceLoadRowAfterRow(const float4* data, int pitch, int width, int height, float4* out)
8  {
9      float4 sum = make_float4(0,0,0,0);
10     int xin = blockIdx.x*blockDim.x + threadIdx.x;
11     int yin = blockIdx.y*blockDim.y + threadIdx.y;
12     for(int y = yin; y < height; y+=blockDim.y)
13     {
14         float4 in = data[xin + y*pitch];
15         sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
16     }
17     out[xin + yin*gridDim.x*blockDim.x] = sum;
18 }
32 __global__ void texLoadRowAfterRow( int width, int height, float4* out)
33 {
34     float4 sum = make_float4(0,0,0,0);
35     int xin = blockIdx.x*blockDim.x + threadIdx.x;
36     int yin = blockIdx.y*blockDim.y + threadIdx.y;
37     for(int y = yin; y < height; y+=blockDim.y)
38     {
39         float4 in = tex2D(myTex,xin,y);
40         sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
41     }
42     out[xin + yin*gridDim.x*blockDim.x] = sum;
43 }
```

RowAfterRow: Texture vs Global



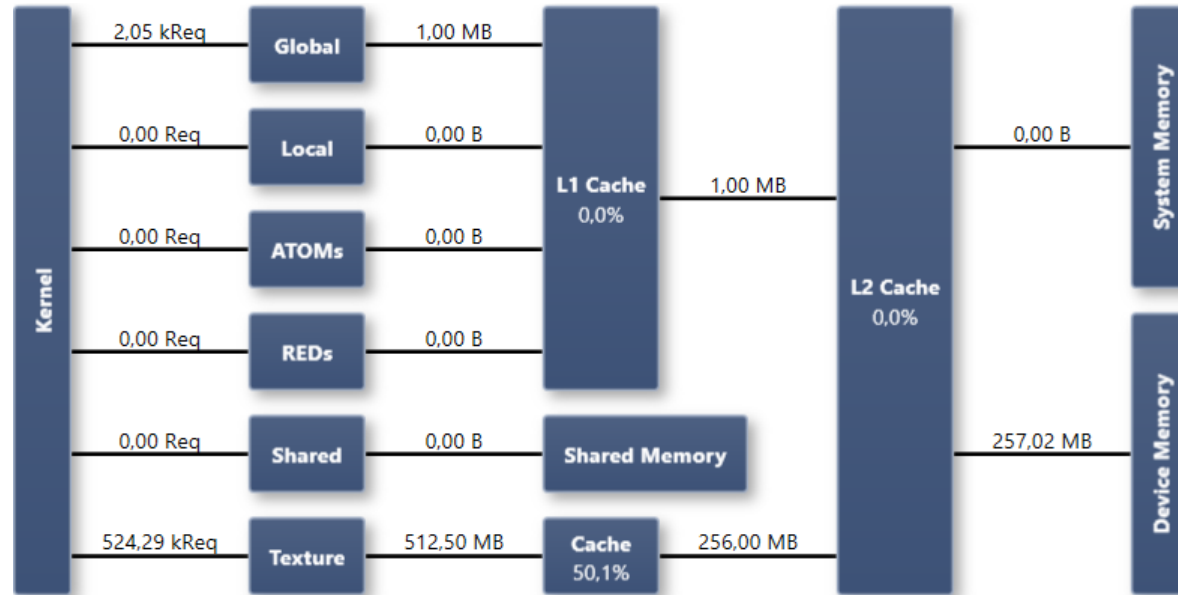
global linear (cudaMalloc) 177GB/s

RowAfterRow: Texture vs Global



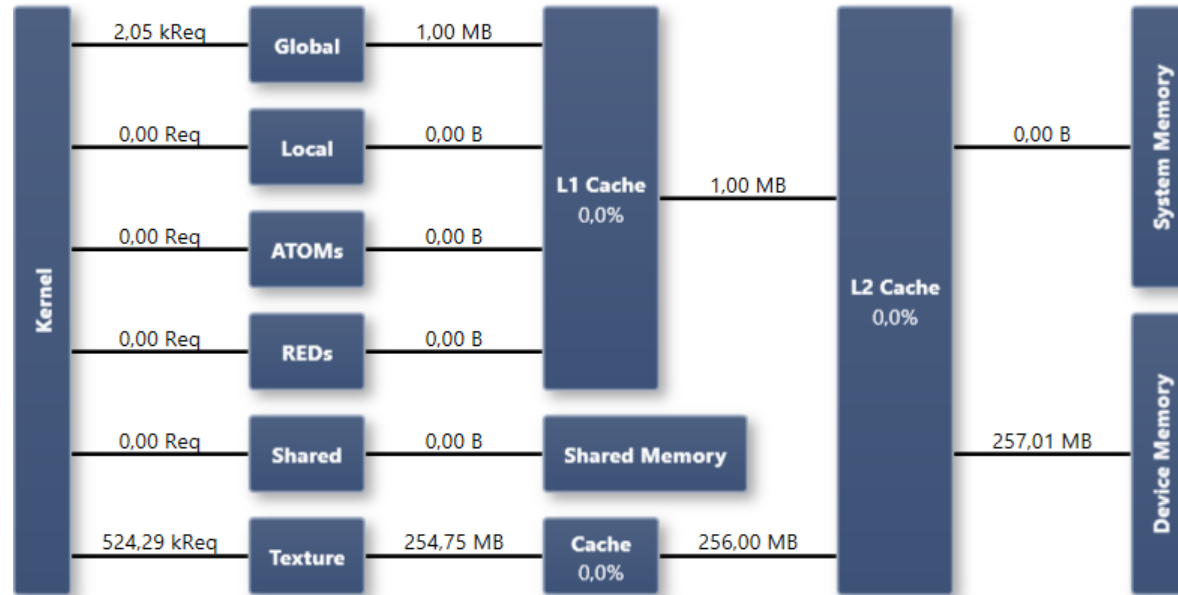
global 2D (cudaMallocPitch) 178GB/s

RowAfterRow: Texture vs Global



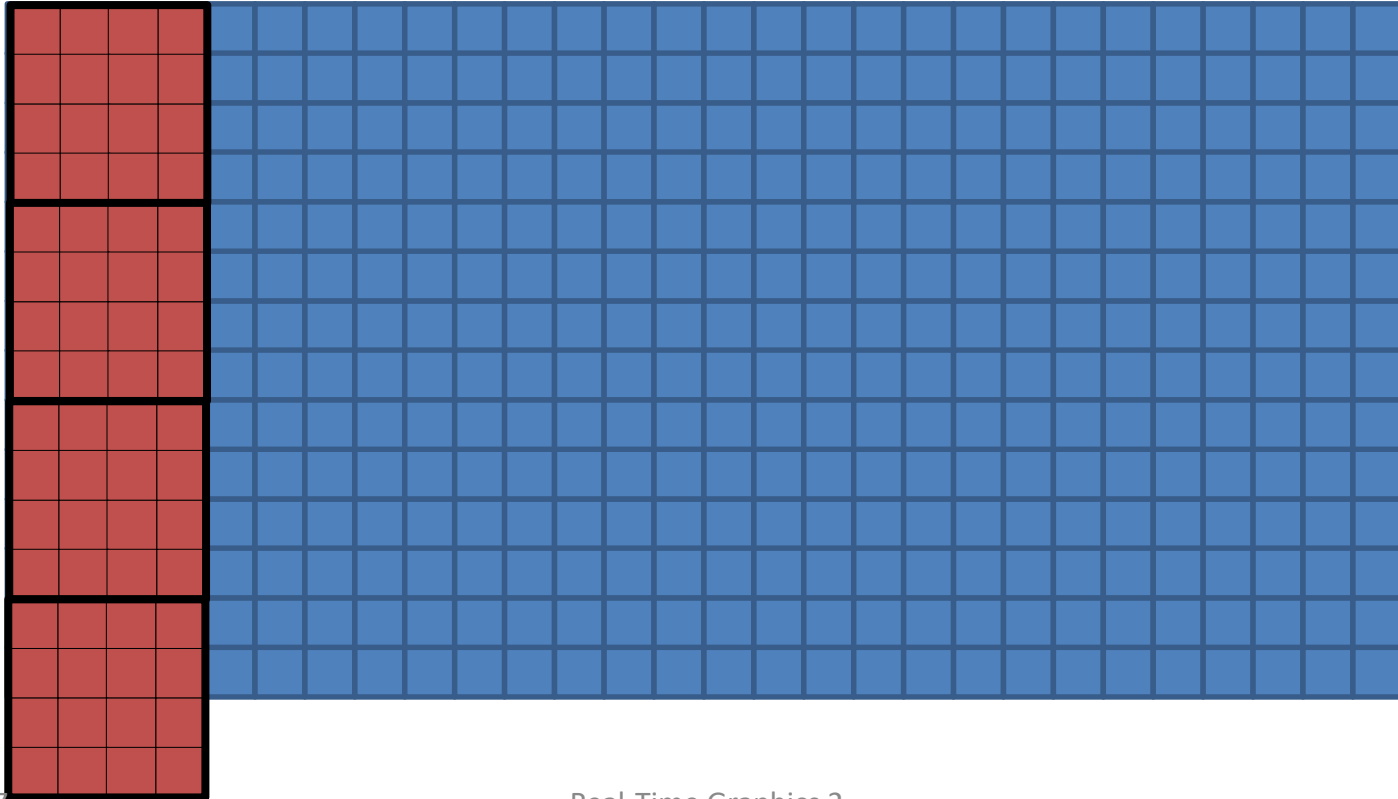
texture (cudaArray) 174GB/s

RowAfterRow: Texture vs Global



texture (cudaMallocPitch) 175GB/s

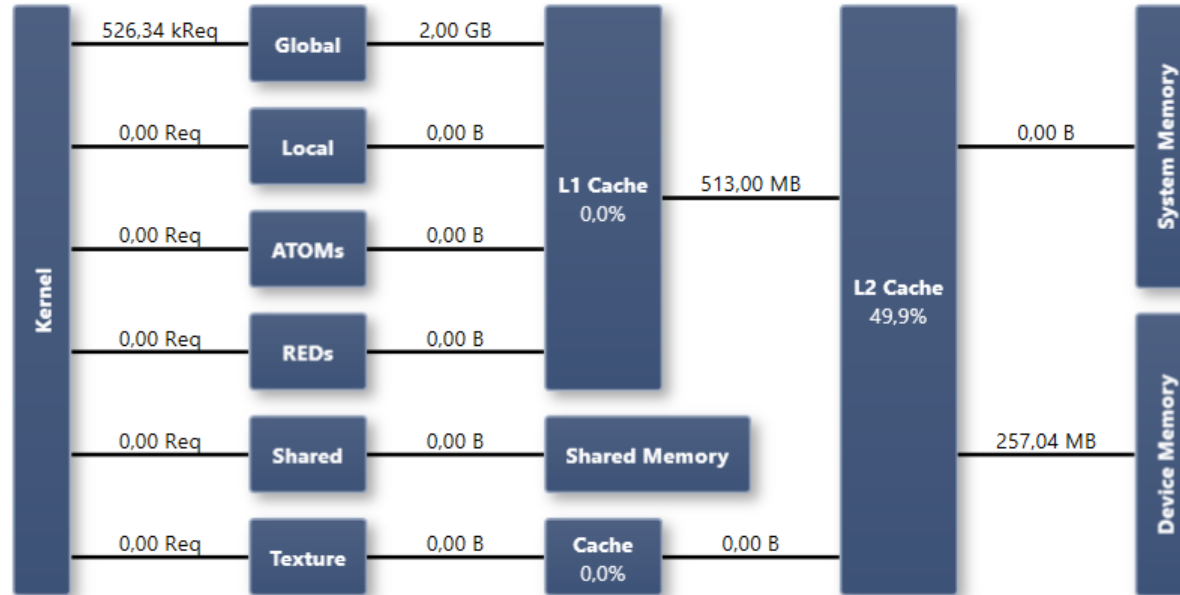
ColumnAfterColumn



ColumnAfterColumn: Texture vs Global

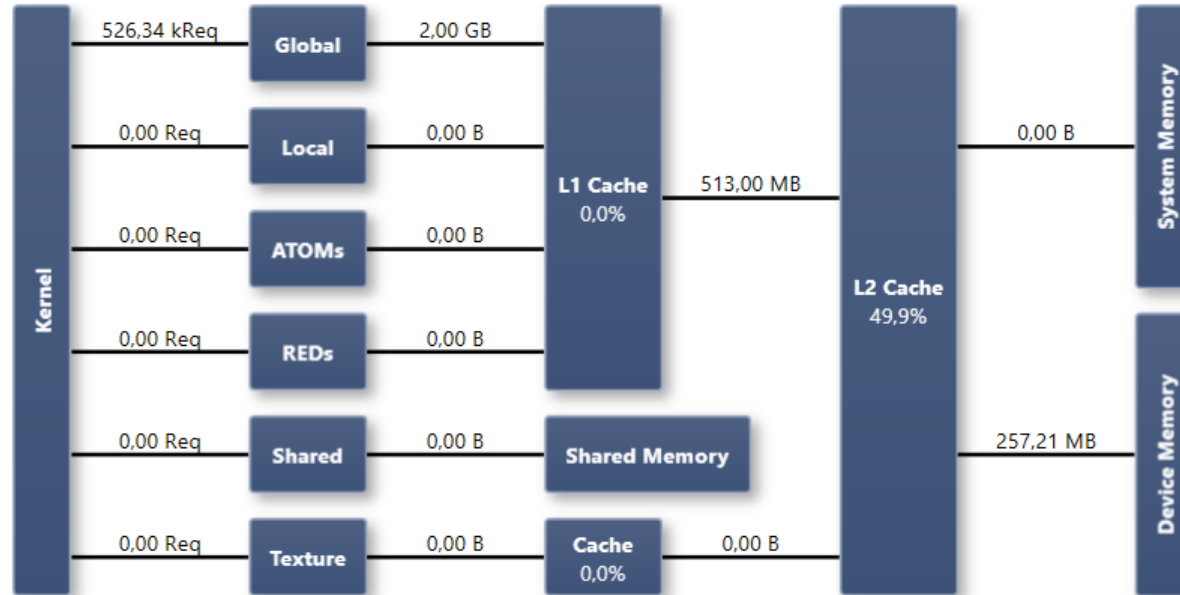
```
19 □ __global__ void deviceLoadColumnAfterColumn(const float4* data, int pitch, int width, int height, float4* out)
20 {
21     float4 sum = make_float4(0,0,0,0);
22     int xin = blockIdx.x*blockDim.x + threadIdx.x;
23     int yin = blockIdx.y*blockDim.y + threadIdx.y;
24     for(int x = xin; x < width; x+=blockDim.x)
25     {
26         float4 in = data[yin + x*pitch];
27         sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
28     }
29     out[xin + yin*gridDim.x*blockDim.x] = sum;
30 }
44 □ __global__ void texLoadColumnAfterColumn( int width, int height, float4* out)
45 {
46     float4 sum = make_float4(0,0,0,0);
47     int xin = blockIdx.x*blockDim.x + threadIdx.x;
48     int yin = blockIdx.y*blockDim.y + threadIdx.y;
49     for(int x = xin; x < width; x+=blockDim.x)
50     {
51         float4 in = tex2D(myTex,yin,x);
52         sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
53     }
54     out[xin + yin*gridDim.x*blockDim.x] = sum;
55 }
```

ColumnAfterColumn: Texture vs Global



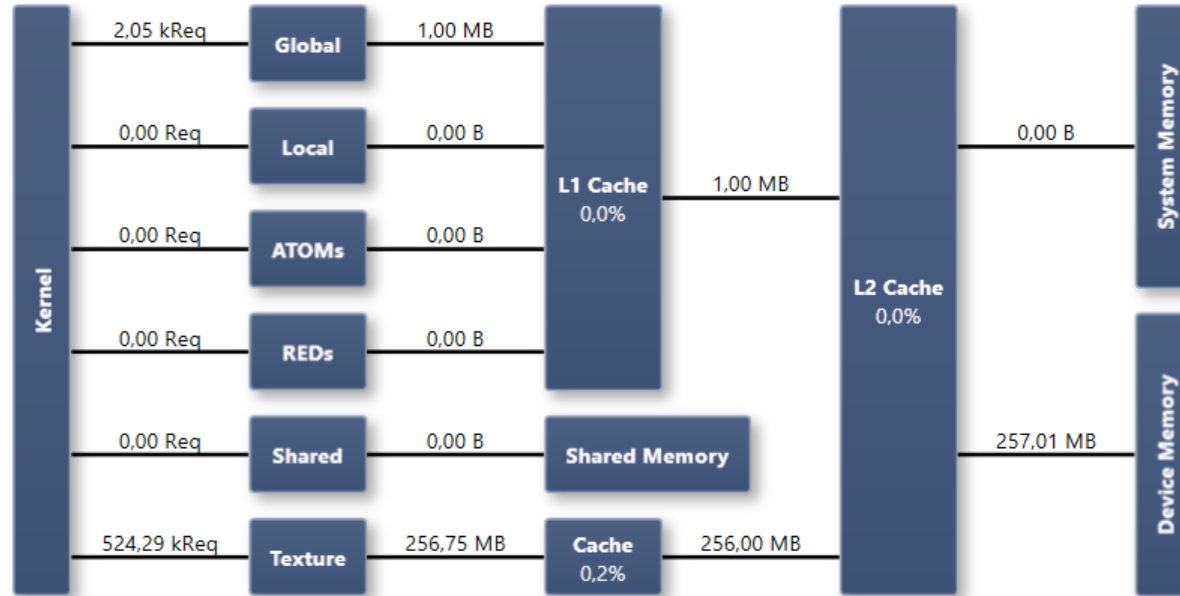
global linear (cudaMalloc) 46GB/s

ColumnAfterColumn: Texture vs Global



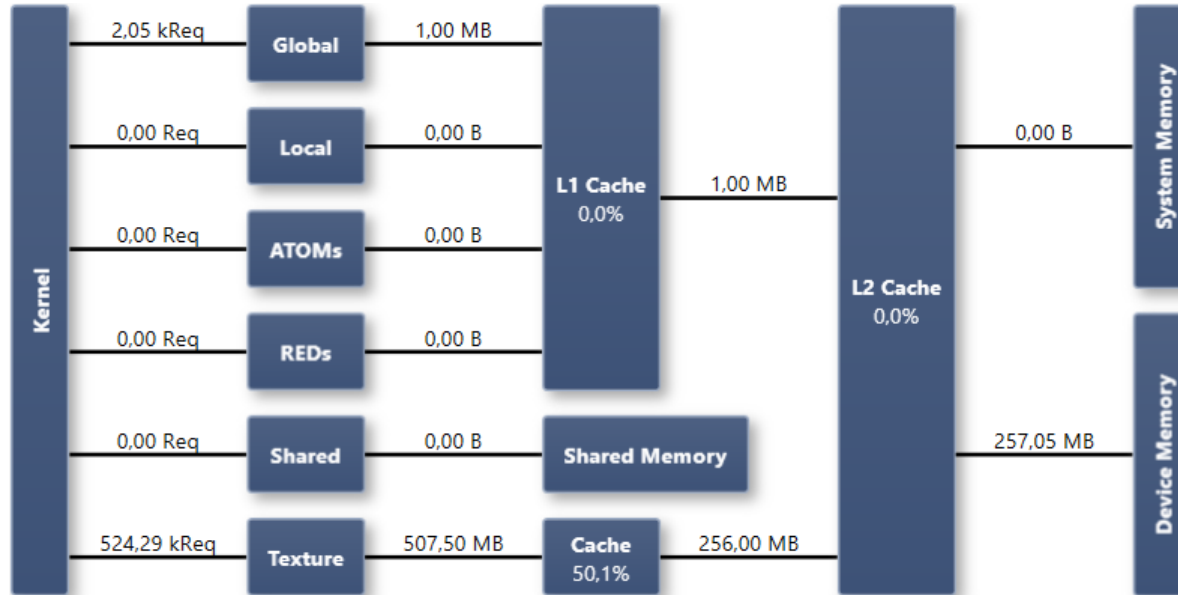
global 2D (cudaMallocPitch) 46GB/s

ColumnAfterColumn: Texture vs Global



texture (cudaArray) 171GB/s

ColumnAfterColumn: Texture vs Global



texture (cudaMallocPitch) 63GB/s

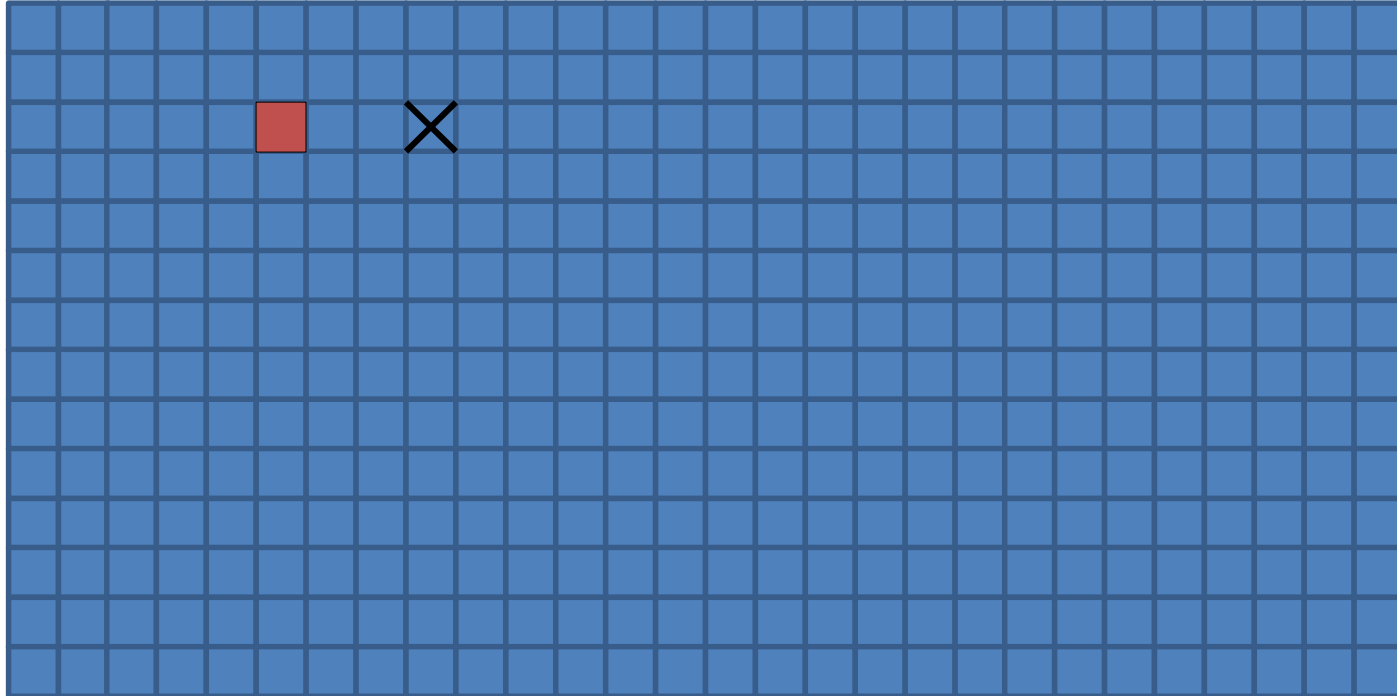
Texture vs Global: Sliding

globalLinRowAfterRow	1.50963ms	-> 177.815GB/s
global2DRowAfterRow	1.50355ms	-> 178.534GB/s
texArrayRowAfterRow	1.53712ms	-> 174.635GB/s
tex2DMemRowAfterRow	1.53299ms	-> 175.106GB/s

globalLinColumnAfterColumn	5.85152ms	-> 45.8745GB/s
global2DColumnAfterColumn	5.85245ms	-> 45.8672GB/s
texArrayColumnAfterColumn	1.56563ms	-> 171.455GB/s
tex2DMemColumnAfterColumn	4.25344ms	-> 63.1102GB/s

- Row access is similarly efficient with all types
- Column access via cudaArray allows for direct load while others need cache
- Texture access more efficient than L2 cache

Filter X



FilterX: Texture vs Global

```

58 __global__ void globalFilterX(const float4* data, int pitch, int width, int height, float4* out, int offset)
59 {
60     float4 sum = make_float4(0,0,0,0);
61     int xin = blockIdx.x*blockDim.x + threadIdx.x;
62     int yin = blockIdx.y*blockDim.y + threadIdx.y;
63
64     if(xin >= offset && xin < width-offset-1)
65     {
66         for(int x = xin-offset; x <= xin+offset; ++x)
67         {
68             float4 in = data[x + yin*pitch];
69             sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
70         }
71     }
72     yin = yin % blockDim.y;
73     out[xin + yin*gridDim.x*blockDim.x] = sum;
74 }

```

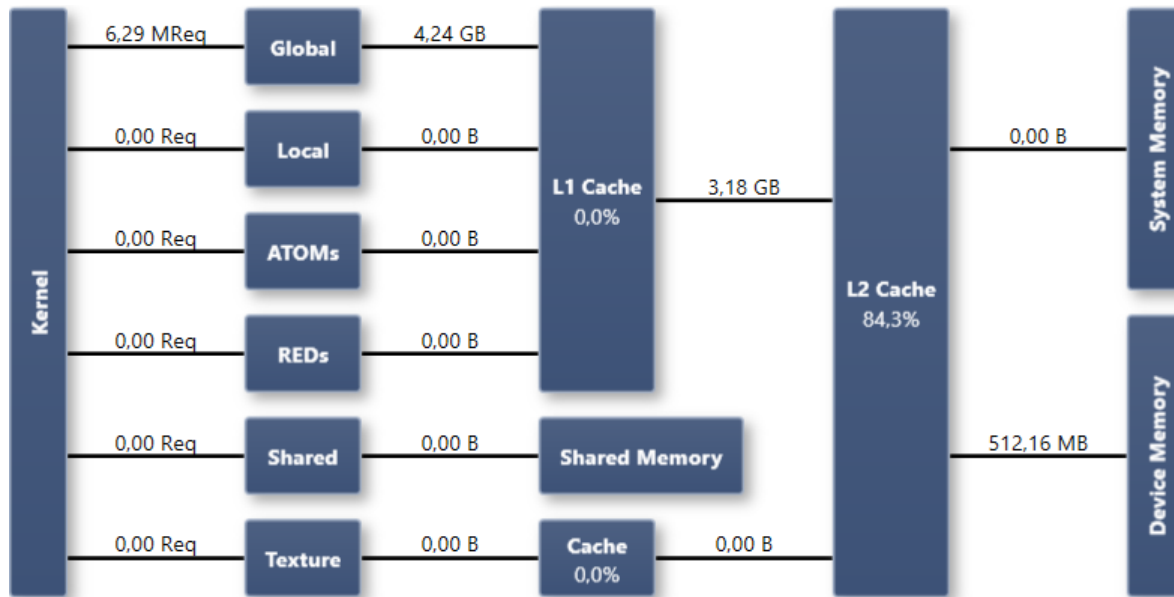
FilterX: Texture vs Global

```

94  __global__ void texFilterX(int width, int height, float4* out, int offset)
95  {
96      float4 sum = make_float4(0,0,0,0);
97      int xin = blockIdx.x*blockDim.x + threadIdx.x;
98      int yin = blockIdx.y*blockDim.y + threadIdx.y;
99
100     for(int x = xin-offset; x <= xin+offset; ++x)
101     {
102         float4 in = tex2D(myTex,x,yin);
103         sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
104     }
105
106     yin = yin % blockDim.y;
107     out[xin + yin*gridDim.x*blockDim.x] = sum;
108 }

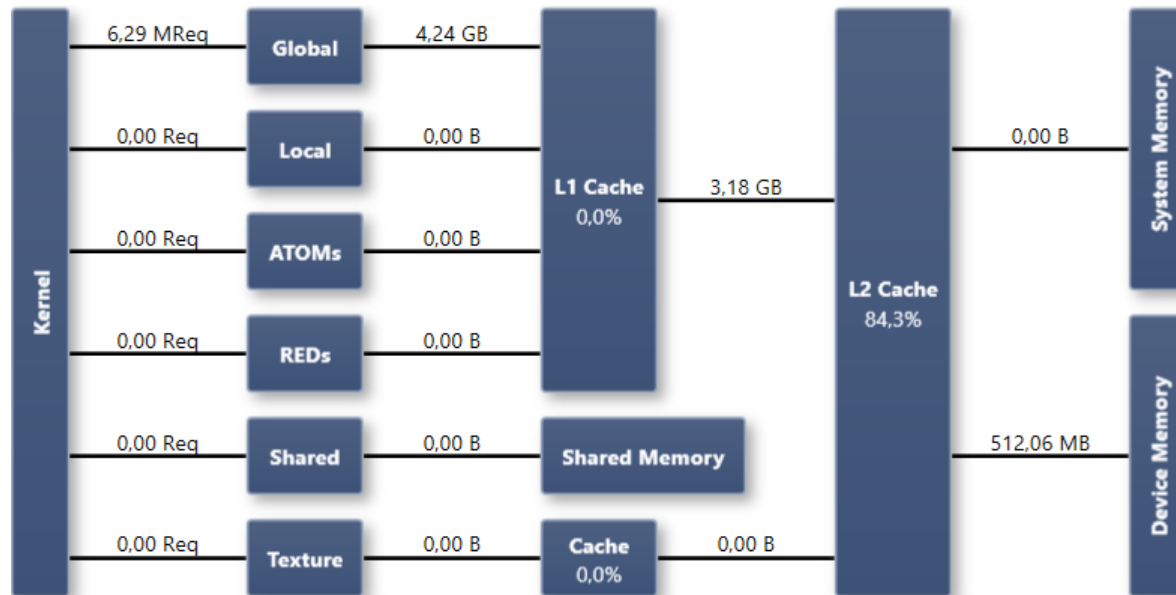
```

FilterX: Texture vs Global



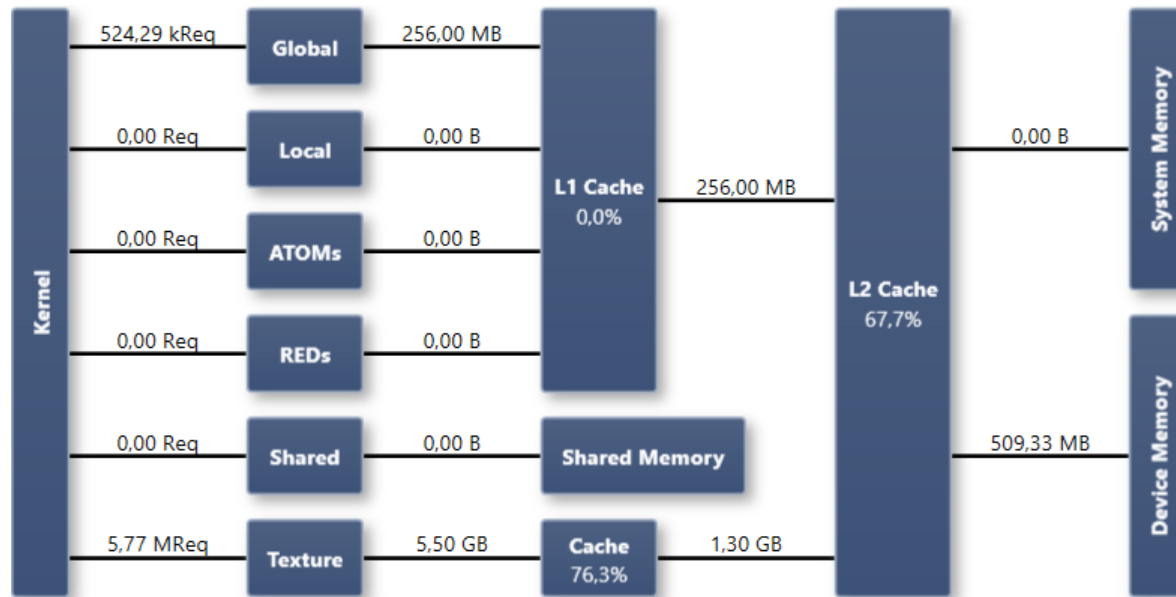
global linear (cudaMalloc) 273GB/s

FilterX: Texture vs Global



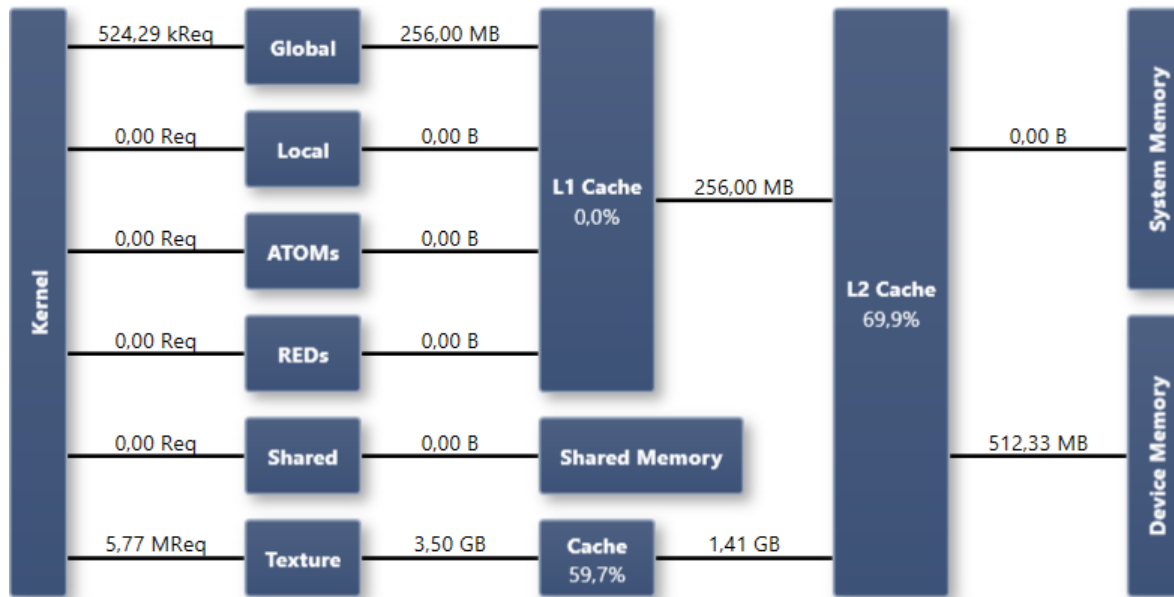
global 2D (cudaMallocPitch) 274GB/s

FilterX: Texture vs Global



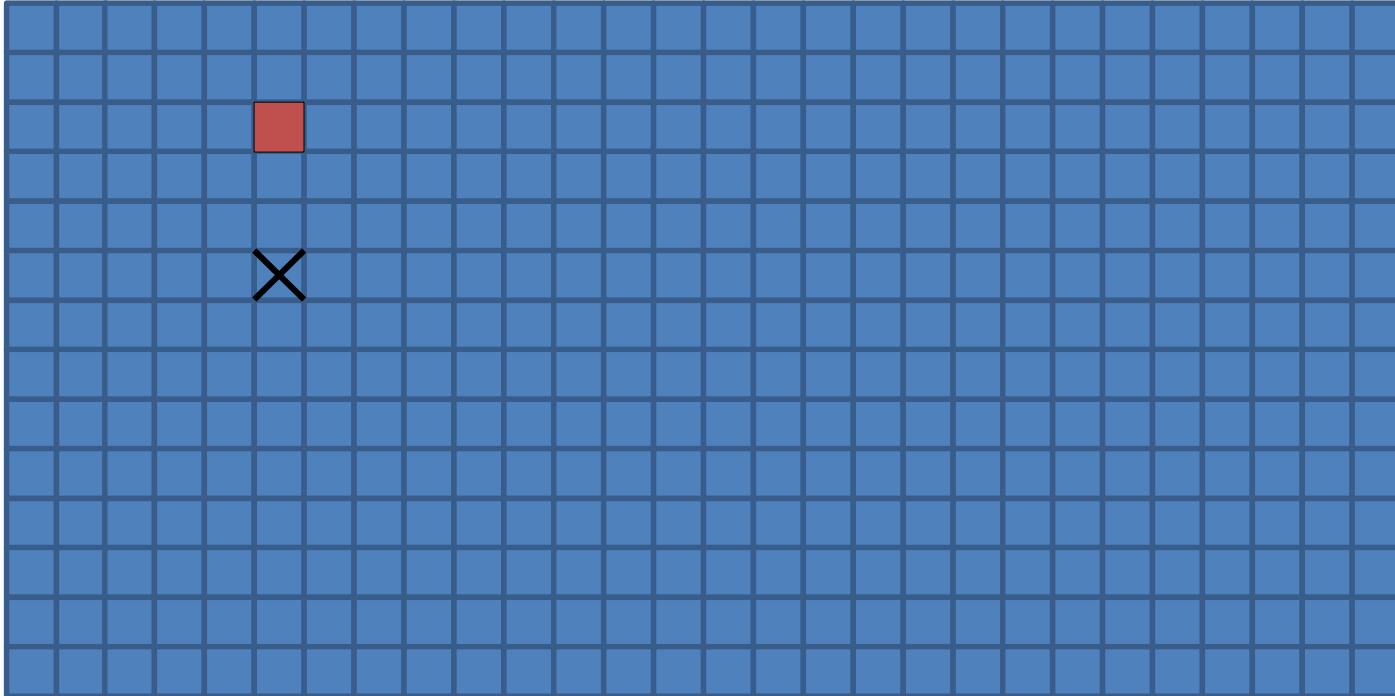
texture (cudaArray) 587GB/s

FilterX: Texture vs Global



texture (cudaMallocPitch) 514GB/s

Filter Y



FilterY: Texture vs Global

```

75 __global__ void globalFilterY(const float4* data, int pitch, int width, int height, float4* out, int offset)
76 {
77     float4 sum = make_float4(0,0,0,0);
78     int xin = blockIdx.x*blockDim.x + threadIdx.x;
79     int yin = blockIdx.y*blockDim.y + threadIdx.y;
80
81     if(yin >= offset && yin < height-offset-1)
82     {
83         for(int y = yin-offset; y <= yin+offset; ++y)
84         {
85             float4 in = data[xin + y*pitch];
86             sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
87         }
88     }
89     yin = yin % blockDim.y;
90     out[xin + yin*gridDim.x*blockDim.x] = sum;
91 }

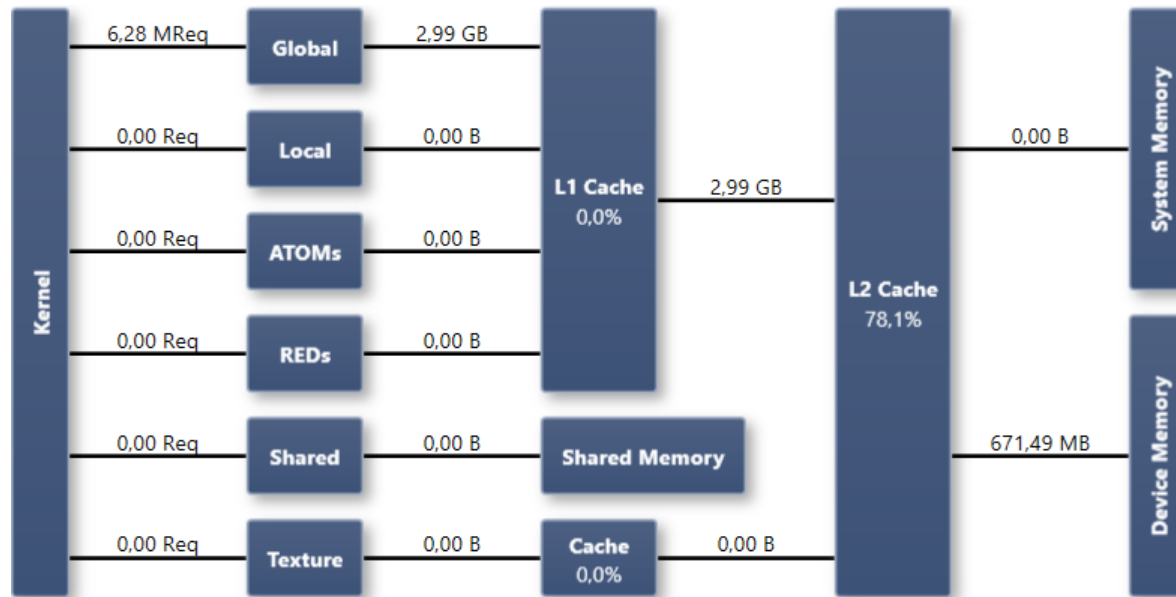
```


FilterY: Texture vs Global

```

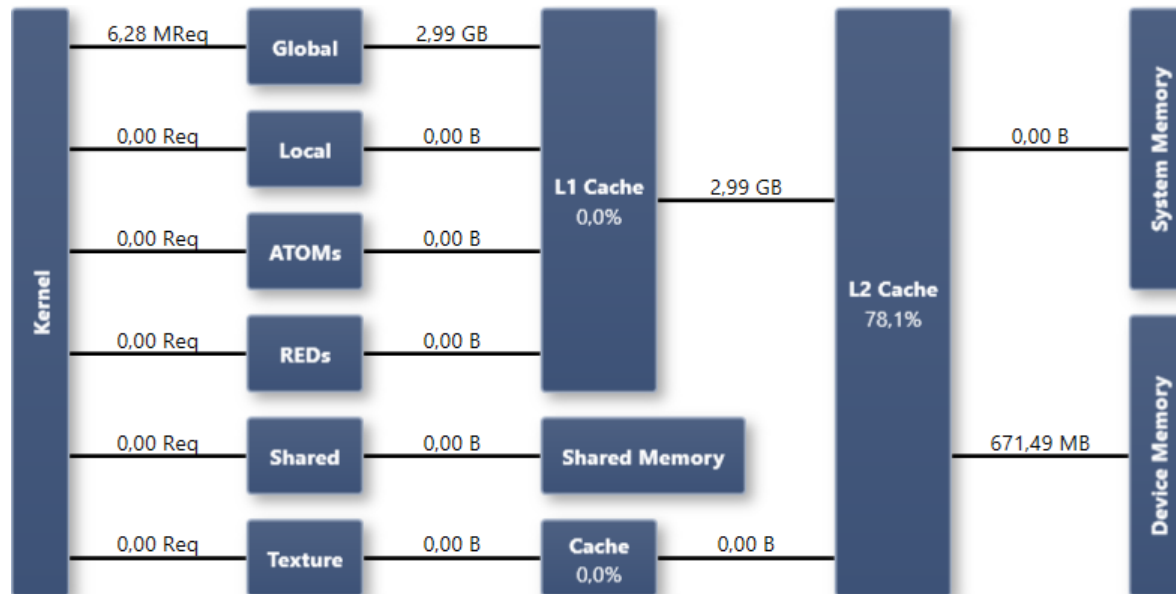
111 ☐ __global__ void texFilterY(int width, int height, float4* out, int offset)
112 {
113     float4 sum = make_float4(0,0,0,0);
114     int xin = blockIdx.x*blockDim.x + threadIdx.x;
115     int yin = blockIdx.y*blockDim.y + threadIdx.y;
116
117     for(int y = yin-offset; y <= yin+offset; ++y)
118     {
119         float4 in = tex2D(myTex,xin,y);
120         sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
121     }
122     yin = yin % blockDim.y;
123     out[xin + yin*gridDim.x*blockDim.x] = sum;
124 }
    
```

FilterY: Texture vs Global



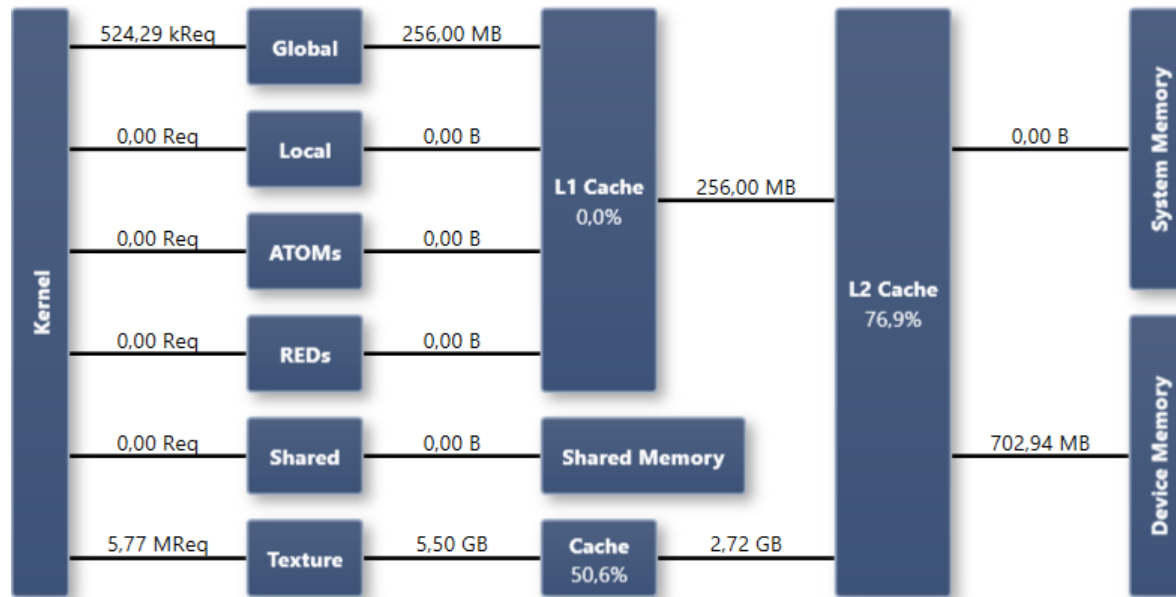
global linear (cudaMalloc) 291GB/s

FilterY: Texture vs Global



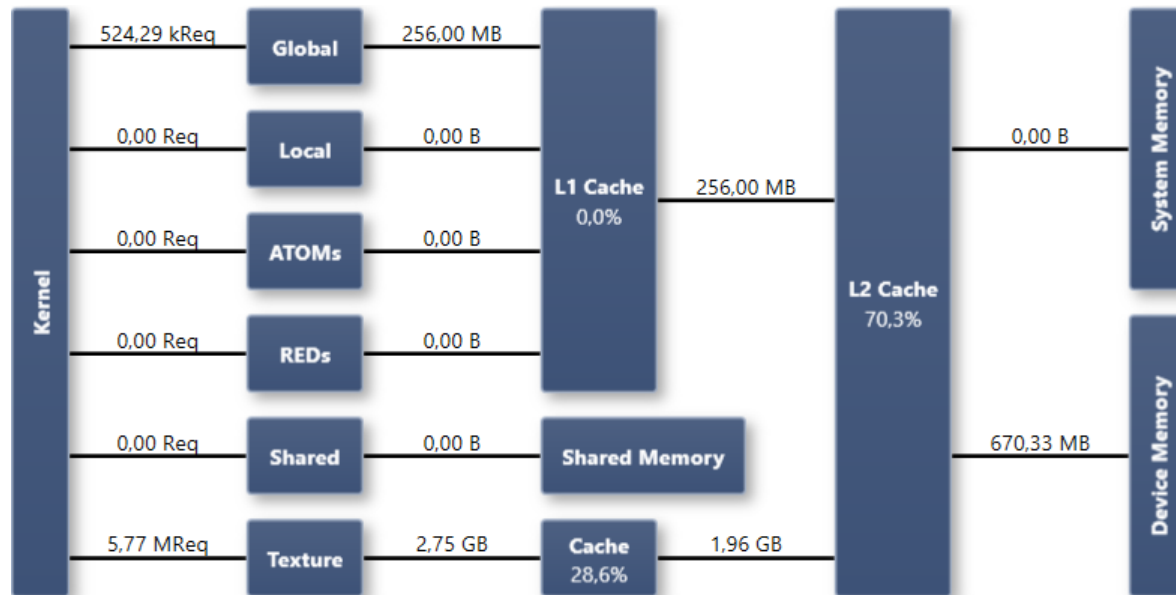
global 2D (cudaMallocPitch) 291GB/s

FilterY: Texture vs Global



texture (cudaArray) 364GB/s

FilterY: Texture vs Global



texture (cudaMallocPitch) 495GB/s

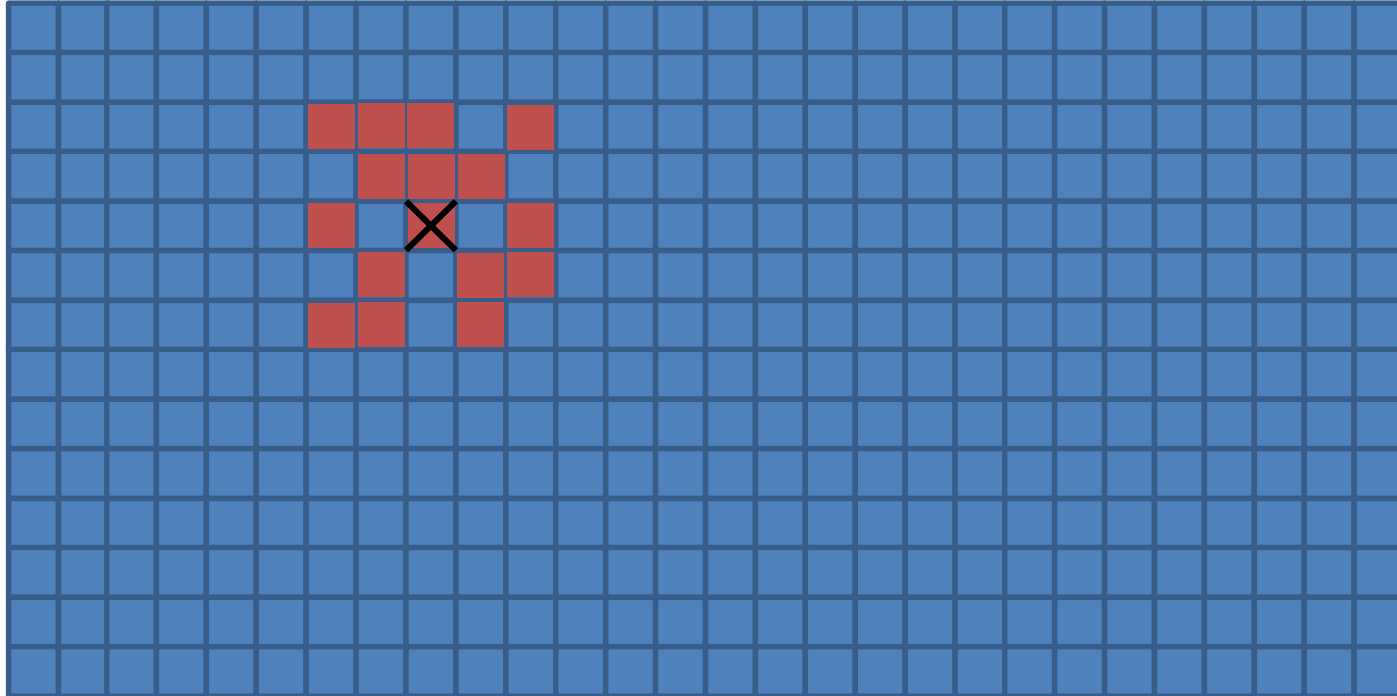
Texture vs Global: Filtering

globalLinFilterX	10.7816ms -> 273.872GB/s
global2DFilterX	10.7704ms -> 274.157GB/s
texArrayFilterX	5.02566ms -> 587.542GB/s
tex2DFilterX	5.73747ms -> 514.65GB/s

globalLinFilterY	10.1444ms -> 291.076GB/s
global2DFilterY	10.1401ms -> 291.2GB/s
texArrayFilterY	8.09706ms -> 364.674GB/s
tex2DFilterY	5.95472ms -> 495.874GB/s

- L2 cache important for all examples,
- Texture cache boosts speed
- Filtering along X is faster than Y

Random Access in vicinity



Random Sampling: Texture vs Global

```
126  template<int Area>
127  □ __global__ void globalRandomSample(const float4* data, int pitch, int width, int height, float4* out, int samples)
128  {
129      float4 sum = make_float4(0,0,0,0);
130      int xin = blockIdx.x*blockDim.x + Area*(threadIdx.x/Area);
131      int yin = blockIdx.y*blockDim.y + Area*(threadIdx.y/Area);
132
133      unsigned int xseed = threadIdx.x *9182 + threadIdx.y*91882
134          + threadIdx.x*threadIdx.y*811 + 72923181;
135      unsigned int yseed = threadIdx.x *981 + threadIdx.y*124523
136          + threadIdx.x*threadIdx.y*327 + 98721121;
137
138      for(int sample = 0; sample < samples; ++sample)
139      {
140          unsigned int x = xseed%Area;
141          unsigned int y = yseed%Area;
142          xseed = (xseed * 1587);
143          yseed = (yseed * 6971);
144          float4 in = data[xin + x + (yin + y)*pitch];
145          sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
146      }
147
148      yin = (yin + threadIdx.y%Area) % blockDim.y;
149      out[xin + threadIdx.x%Area + yin*width] = sum;
150  }
```

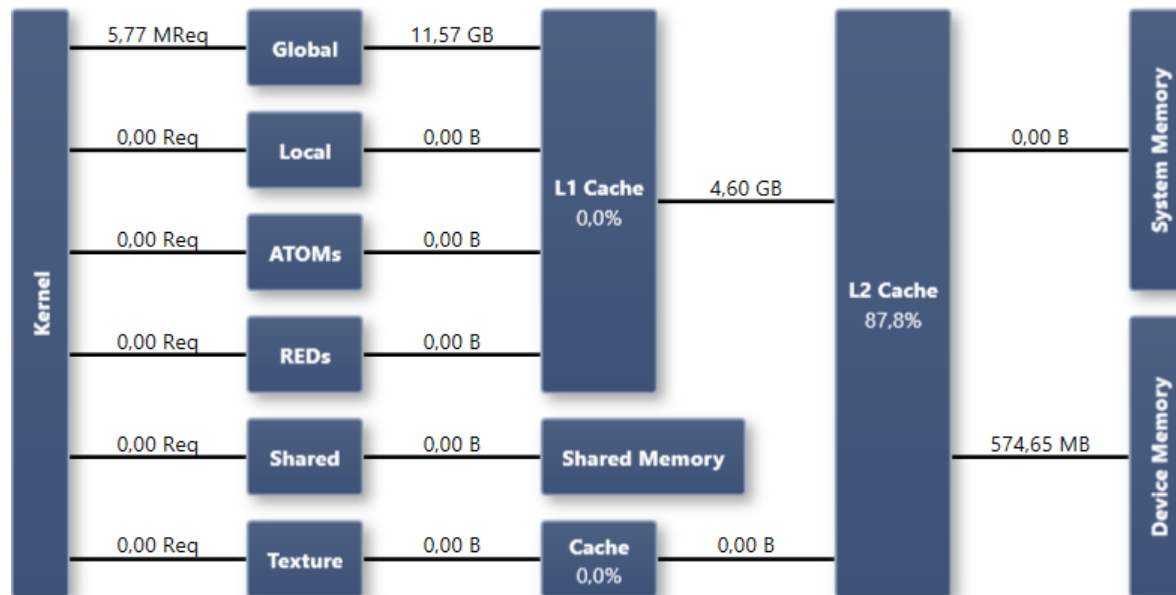

Random Sampling: Texture vs Global

```

152  template<int Area>
153  □ __global__ void texRandomSmple(int width, int height, float4* out, int samples)
154  {
155      float4 sum = make_float4(0,0,0,0);
156      int xin = blockIdx.x*blockDim.x + Area*(threadIdx.x/Area);
157      int yin = blockIdx.y*blockDim.y + Area*(threadIdx.y/Area);
158
159      unsigned int xseed = threadIdx.x * 9182 + threadIdx.y*9182
160          + threadIdx.x*threadIdx.y*811 + 72923181;
161      unsigned int yseed = threadIdx.x * 981 + threadIdx.y*124523
162          + threadIdx.x*threadIdx.y*327 + 98721121;
163
164      for(int sample = 0; sample < samples; ++sample)
165      {
166          unsigned int x = xseed%Area;
167          unsigned int y = yseed%Area;
168          xseed = (xseed * 1587);
169          yseed = (yseed * 6971);
170          float4 in = tex2D(myTex, xin + x, yin + y);
171          sum.x += in.x; sum.y += in.y; sum.z += in.z; sum.w += in.w;
172      }
173
174      yin = (yin + threadIdx.y%Area) % blockDim.y;
175      out[xin + threadIdx.x%Area + yin*width] = sum;
176  }
177

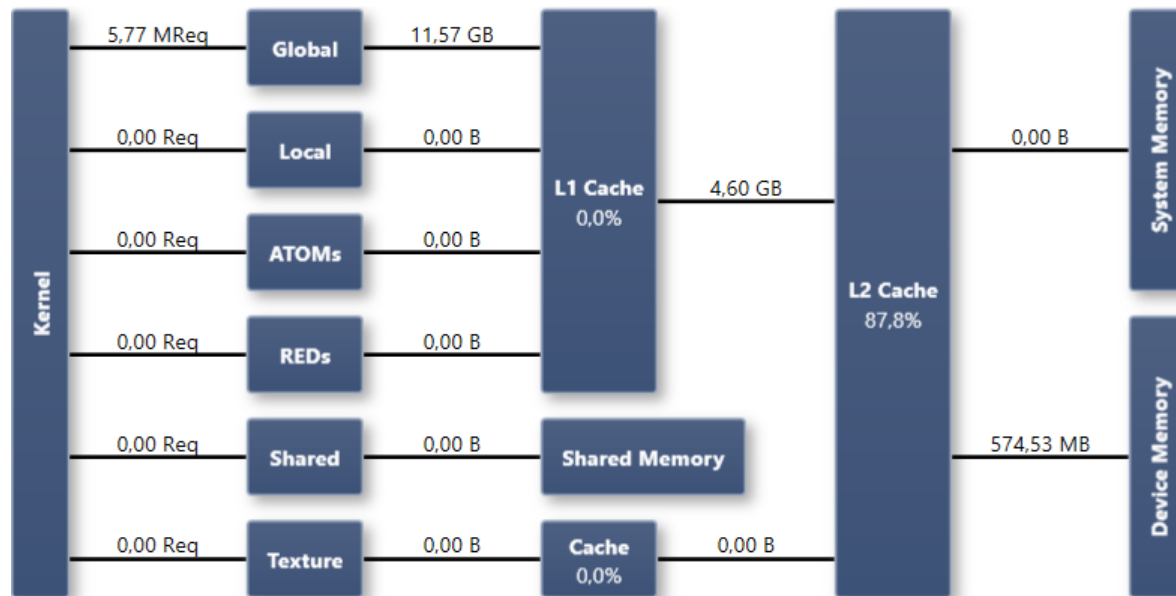
```

Random Sampling: Texture vs Global



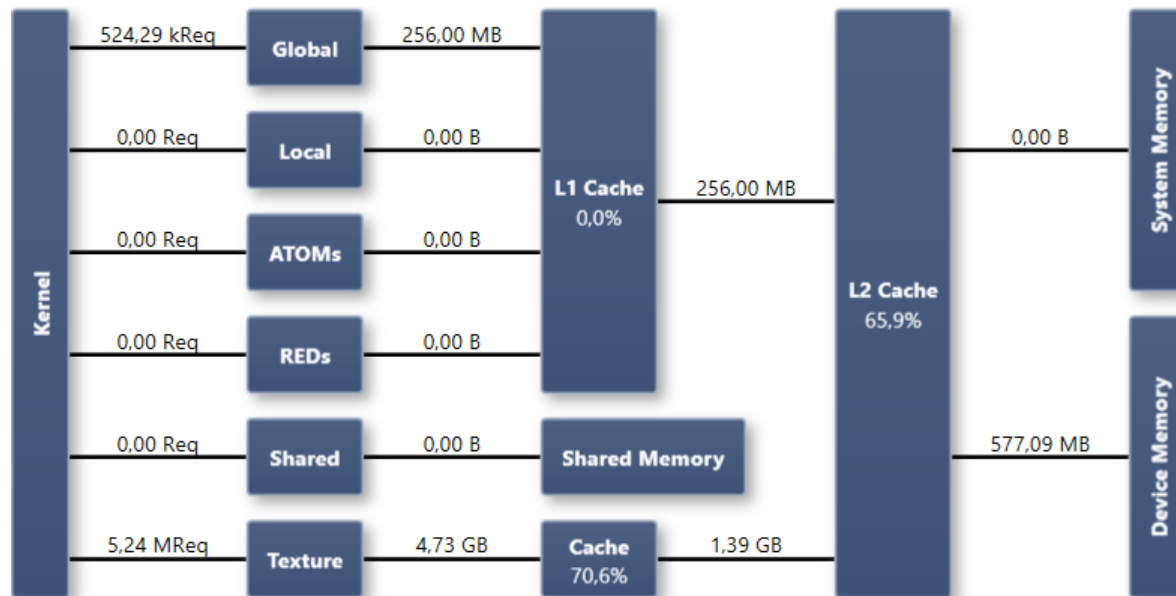
global linear (cudaMalloc) 112GB/s

Random Sampling: Texture vs Global



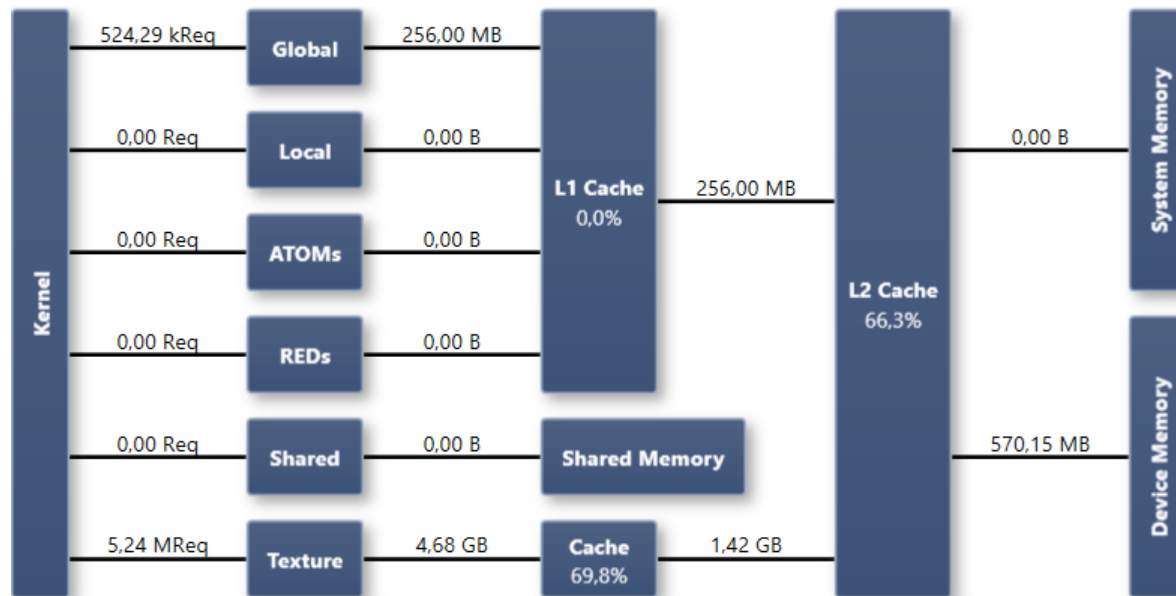
global 2D (cudaMallocPitch) 112GB/s

Random Sampling: Texture vs Global



texture (cudaArray) 355GB/s

Random Sampling: Texture vs Global



texture (cudaMallocPitch) 356GB/s

Texture vs Global: Random Sampling

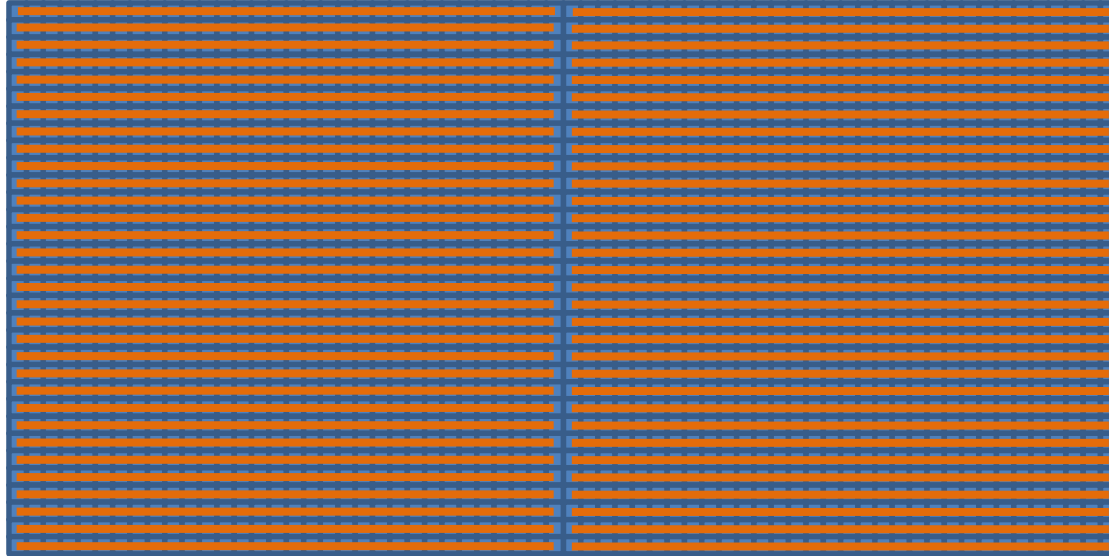
globalLinRandomSample	23.9213ms -> 112.216GB/s
global2DRandomSample	23.9196ms -> 112.224GB/s
texArrayRandomSample	7.54998ms -> 355.544GB/s
tex2DRandomSample	7.53933ms -> 356.047GB/s

- L2 cache important for all examples
- Texture cache boosts speed

Texture vs Global Conclusion

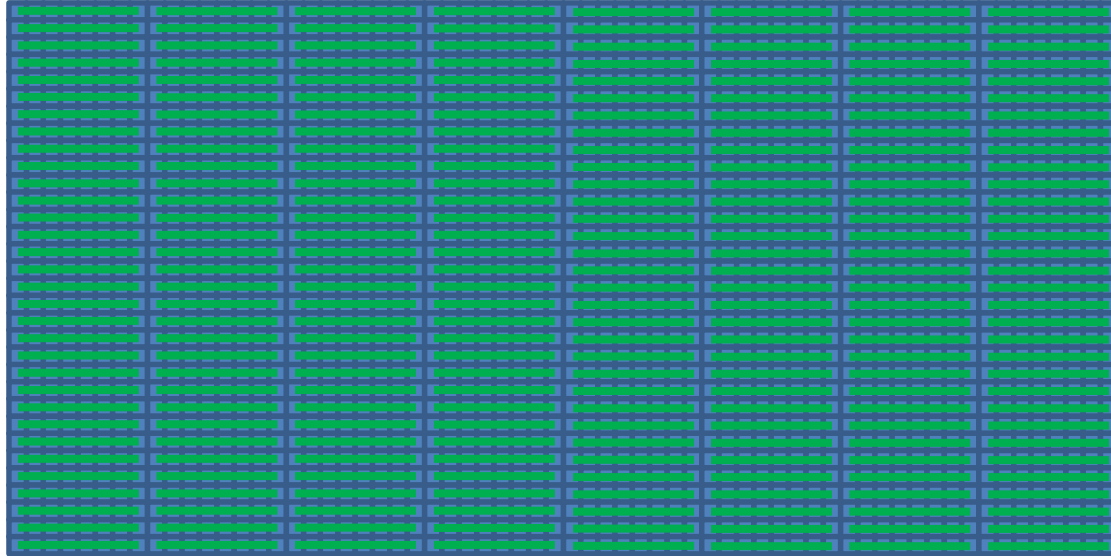
- Textures always perform as good or better
 - With textures performance less dependent on access direction/pattern
 - Textures put less stress on L2 cache
 - Texture cache separate from L1 cache (L1 free for other tasks)
 - Features: border handling, interpolation, conversion

Array vs Pitched Textures



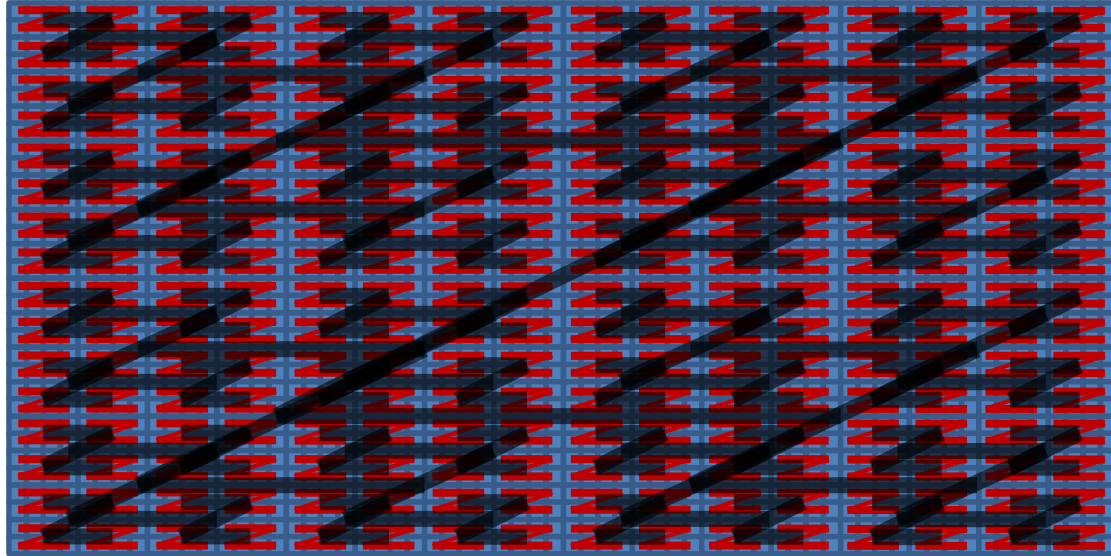
- Pitched global access: cache line size 128 byte

Array vs Pitched Textures



- Pitched global access: cache line size 128 byte
- Pitched texture access: cache line size 32 byte

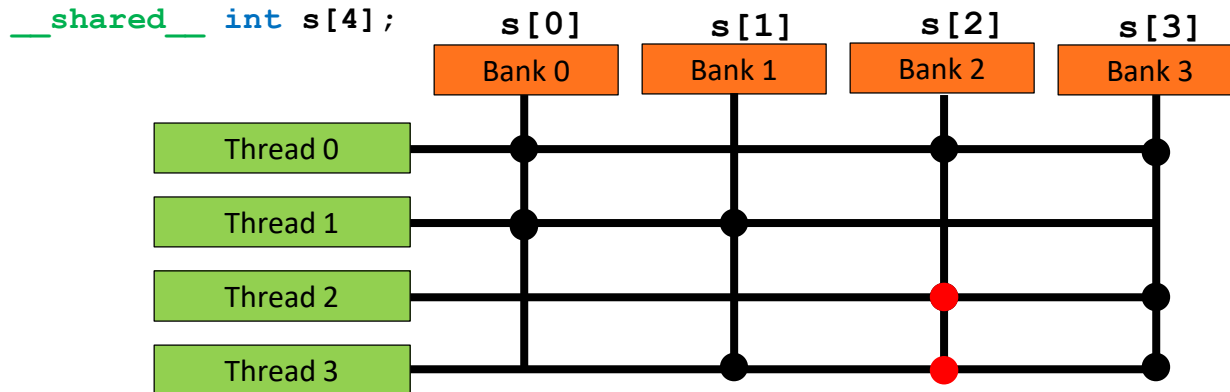
Array vs Pitched Textures



- Pitched global access: cache line size 128 byte
- Pitched texture access: cache line size 32 byte
- Array texture access: space filling curve (maybe Z curve)
+ cache line size 32 byte

Shared memory recap

- Shared access within one block (lifetime: block)
- Located on multiprocessor → very fast
- Limited available size on multiprocessor
- Crossbar: simultaneous access to distinct banks



Shared memory

- Accessed via crossbar → access pattern important
- Different behavior for different architectures
 - Pre-Fermi:
 - access issued per 16 threads
 - 16 banks a 32 bits per 2 clock cycles

$$\text{bank} = (\text{address}/4) \% 16$$
 - Fermi:
 - access issued per 32 threads
 - 32 banks a 32 bits per 2 clock cycles

$$\text{bank} = (\text{address}/4) \% 32$$
 - Kepler:
 - access issued per 32 threads
 - 32 banks a 64 bits per clock cycle
 - 64 bit mode: two threads can access any part of a 64 bit word
 - 32 bit mode: two threads can access any part of a 64 bit word which would fall in the same bank
 - modes can be set using the CUDA API

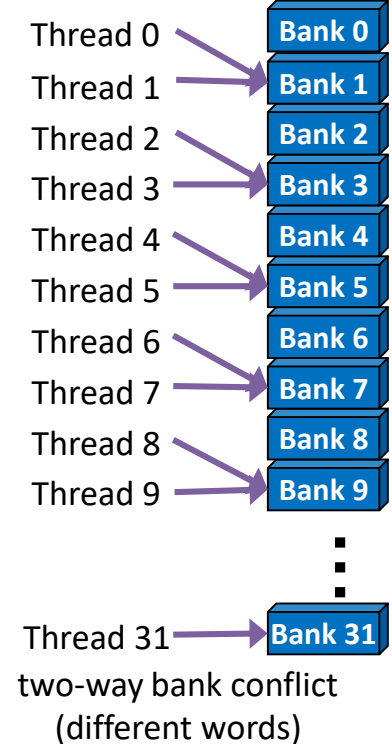
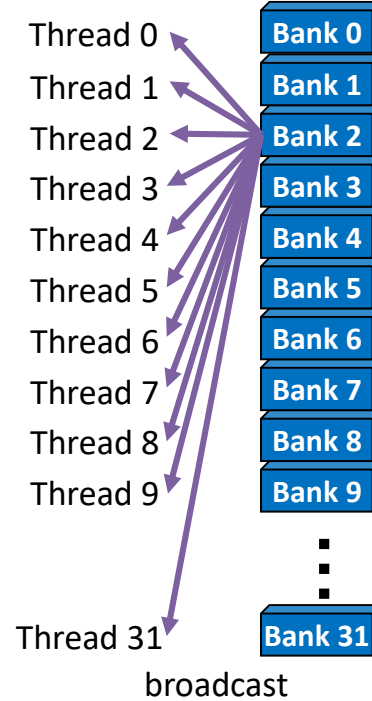
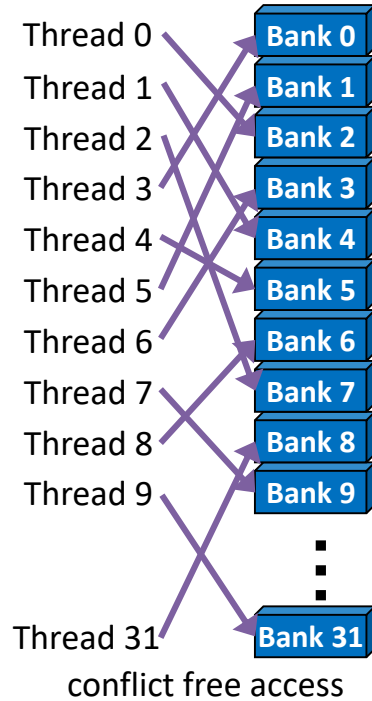
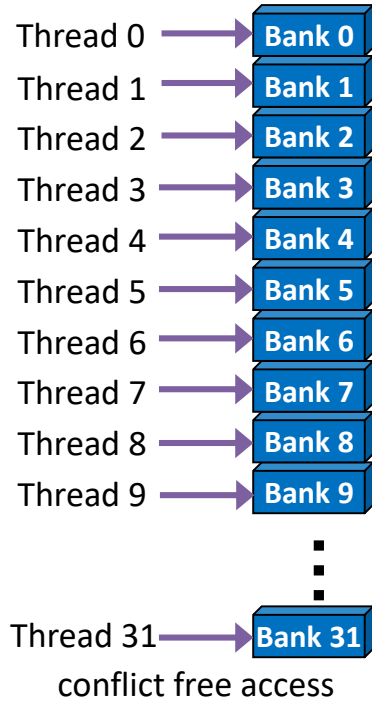
$$\text{bank} = (\text{address}/8) \% 32$$

$$\text{bank} = (\text{address}/4) \% 32$$

Shared memory cont'd

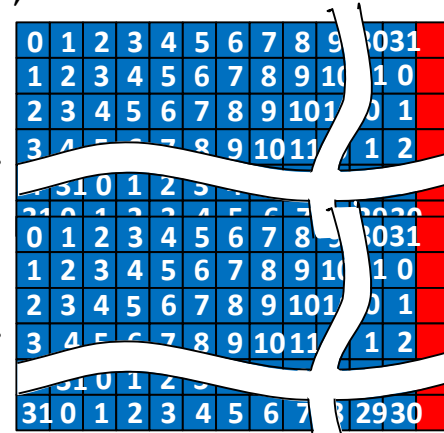
- Key to good performance is the access pattern
- Conflict free access: all threads within a warp access different banks (mind Kepler exceptions)
- Multicast: all threads accessing the same word are served with one transaction
- Serialization: multiple access conflicts will be serialized
- Introduce padding to avoid bank conflicts

Shared memory cont'd



Shared memory cont'd

```
__global__ void kernel(...)
{
    __shared__ float mydata[32*(32 + 1)];
    ...
    float sum = 0;
    for(uint i = 0; i < 32; ++i)
        sum += mydata[threadIdx.x + i*33];
    ...
    sum = 0;
    for(uint i = 0; i < 32; ++i)
        sum += mydata[threadIdx.x*33 + i];
    ...
}
```



0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32
2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33
3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34
4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35
5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36
6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37
7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38
8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39
9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40
10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41
11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42
12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43
13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44
14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47
17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49
19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50
20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51
21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52
22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53
23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54
24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55
25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56
26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57
27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58
28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59
29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60
30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61
31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62

Shared memory use cases

- Inter-thread communication

```
__global__ void kernel(...)
{
    __shared__ bool run;
    run = true;    //for cc>=2 faster than if(threadId == 0)
    __syncthreads();
    while(run)
    {
        ...
        if(found_it)
            run = false;
        __syncthreads();
    }
}
```

Shared memory use cases

- Inter-thread communication
- Reduce global memory access → manual cache

```
__global__ void kernel(float* global_data, ...)
{
    extern __shared__ float data[];
    uint linid = blockIdx.x*blockDim.x + threadIdx.x;
    //load
    data[threadIdx.x] = global_data[linid];
    __syncthreads();
    for(uint it = 0; it < max_it; ++it)
        calc_iteration(data); //calc
    __syncthreads();
    //write back
    global_data[linid] = data[threadIdx.x];
}
```

Shared memory use cases

- Inter-thread communication
- Reduce global memory access → manual cache
- Adjust global memory access pattern

```
__global__ void transp(float* global_data, float* global_data2)
{
    extern __shared__ float data[];
    uint linid1 = blockIdx.x*32 + threadIdx.x +
        blockIdx.y*32*width;
    uint linid2 = blockIdx.x*32*width + threadIdx.x +
        blockIdx.y*32;
    for(uint i = 0; i < 32; ++i)
        data[threadIdx.x + i*33] = global_data[linid1 + i*width];
    __syncthreads();
    for(uint j = 0; j < 32; ++j)
        global_data2[linid2 + j*width] = data[threadIdx.x*33 + j] ;
}
```

Shared memory use cases

- Inter-thread communication
- Reduce global memory access → manual cache
- Adjust global memory access pattern
- Indexed access

```
__global__ void kernel(...)
{
    uint mydata[8]; //will be spilled to local memory
    for(uint i = 0; i < 8; ++i)
        mydata[i] = complexfunc(i, threadIdx.x);
    uint res = 0;
    for(uint i = 0; i < 64; ++i)
        res += secondfunc(mydata[(threadIdx.x + i) % 8],
                           mydata[i*threadIdx.x % 8]);
}
```

Shared memory use cases

- Inter-thread communication
- Reduce global memory access → manual cache
- Adjust global memory access pattern
- Indexed access

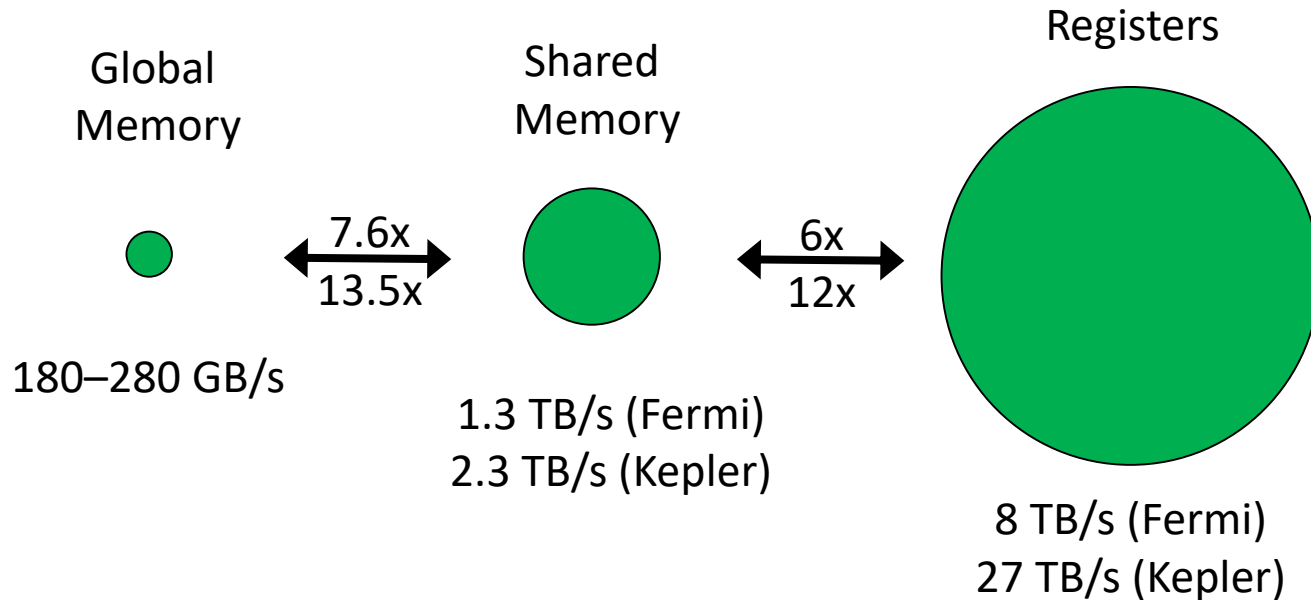
```
__global__ void kernel(...)
{
    __shared__ uint allmydata[8*BlockSize];
    uint *mydata = allmydata + 8*threadIdx.x;
    for(uint i = 0; i < 8; ++i)
        mydata[i] = complexfunc(i, threadIdx.x);
    uint res = 0;
    for(uint i = 0; i < 64; ++i)
        res += secondfunc(mydata[(threadIdx.x + i) % 8],
                           mydata[i*threadIdx.x % 8]);
}
```

Shared memory use cases

- Inter-thread communication
- Reduce global memory access → manual cache
- Adjust global memory access pattern
- Indexed access
- Combine costly operations

```
__global__ void kernel(uint *global_count, ...)
{
    __shared__ uint blockcount;
    blockcount = 0;
    __syncthreads();
    uint myoffset = atomicAdd(&blockcount, myadd);
    __syncthreads();
    if(threadIdx.x == 0)
        blockcount = atomicAdd(global_count, blockcount);
    __syncthreads();
    myoffset += blockcount;
}
```

Registers vs Shared Memory vs Global



Questions

