

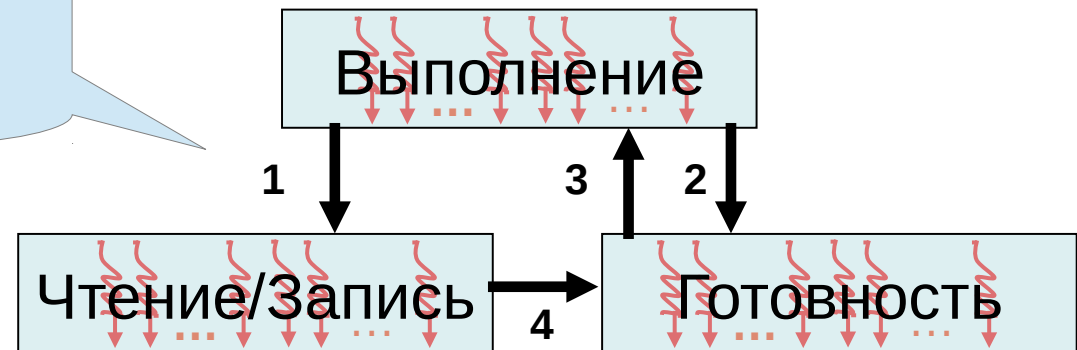
# Лекция 3

## СОДЕРЖАНИЕ

- оптимизация нагрузки GPU;
- иерархия памяти;
- объединение запросов к глобальной памяти (coalescing);

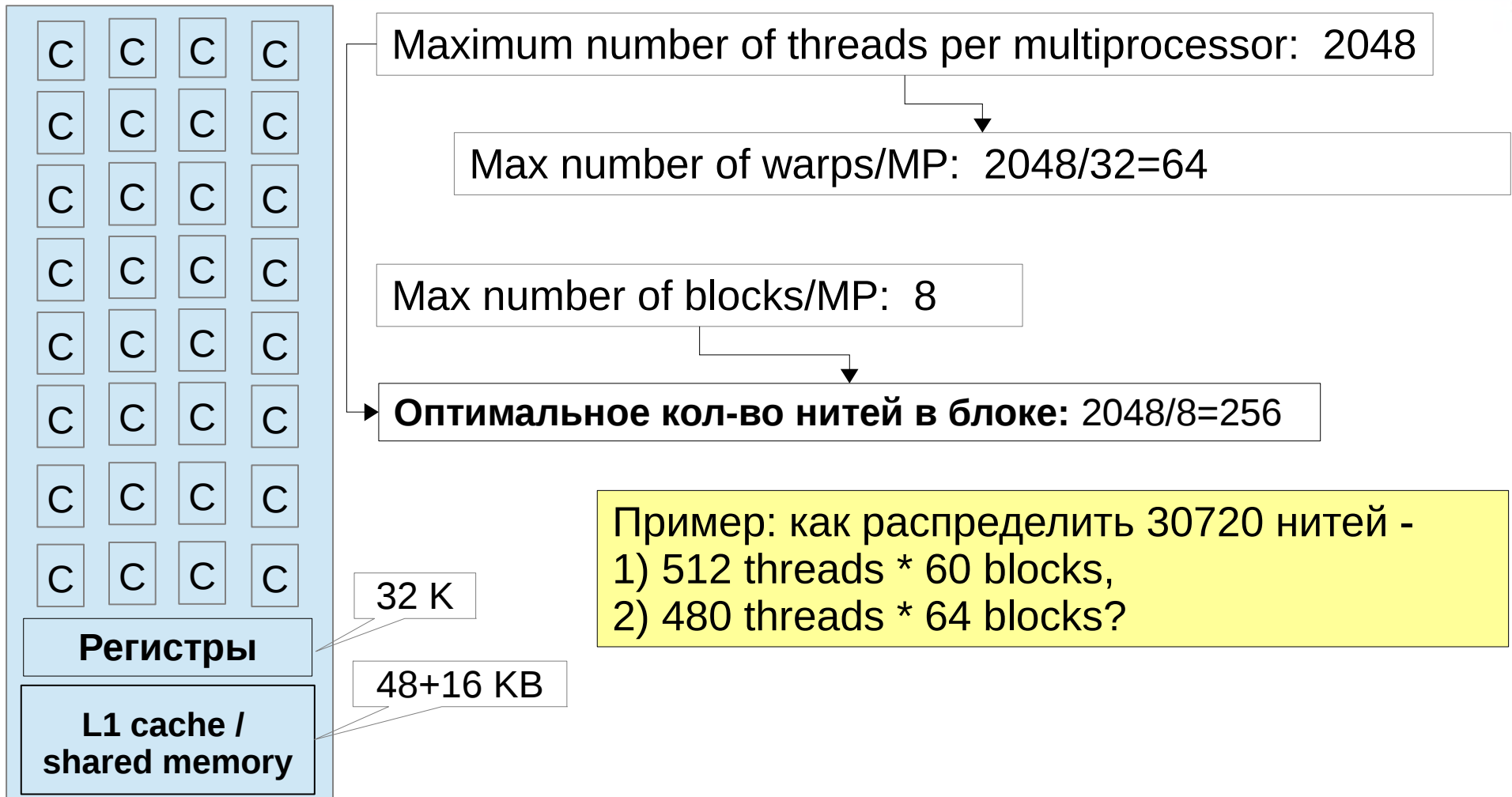
## Оптимальное количество блоков и нитей в блоках (сокращение латентности)

Преимущества  
перекрывания/многозадачности:



Кол-во варпов  $\% 32 == 0$  && Кол-во варпов  $/ 32 > 1$   
Кол-во блоков  $\geq$  Кол-во мультипроцессоров

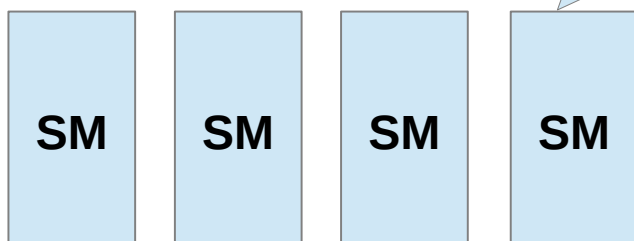
## Оптимальное количество блоков и нитей в блоках (заполняемость мультипроцессоров)



## Оптимальное количество блоков и нитей в блоках (заполняемость мультимикроспроцессоров)

1)

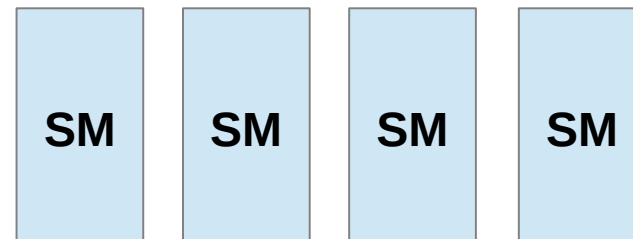
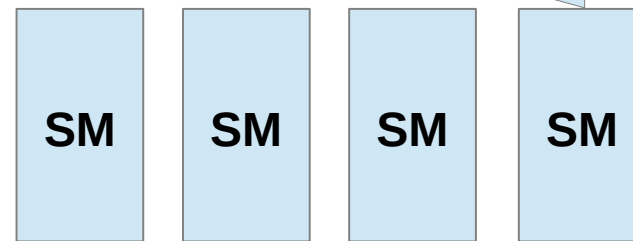
4 blocks \* 16 warps \* 32 threads



1-ая  
загрузка

2)

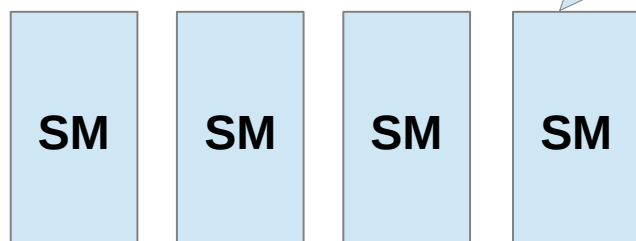
4 blocks \* 15 warps \* 32 threads



## Оптимальное количество блоков и нитей в блоках (заполняемость мультипроцессоров)

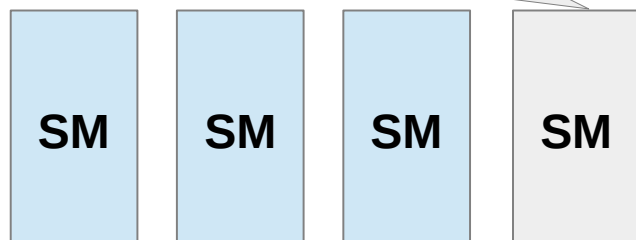
1)

4 blocks \* 16 warps \* 32 threads



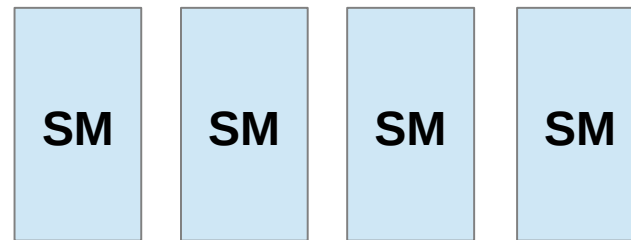
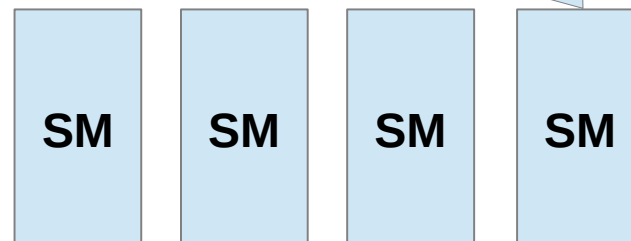
2-ая  
загрузка

Мультипроцессор простаивает



2) сбалансированная нагрузка

4 blocks \* 15 warps \* 32 threads



## Оптимальное количество блоков и нитей в блоках (заполняемость мультипроцессоров)

Для архитектуры Fermi:

- максимальное кол-во варпов, выполняемых одновременно на одном мультипроцессоре равно 48;
- максимальное кол-во блоков выполняемых одновременно на одном мультипроцессоре равно 8;

**Пример.**

32 нити в блоке, 512 блоков в гриде:

Кол-во одновременно выполняемых нитей= $32 \times 8$ ,

Максимальное кол-во= $48 \times 32$ ,

Заполняемость= $32 \times 8 / 48 \times 32 = 1/6 = 16.(6)\%$ .

512 нитей в блоке, 32 блока в гриде:

Кол-во одновременно выполняемых нитей= $512 \times 8$ ,

Максимальное кол-во= $48 \times 32$ ,

Заполняемость= $512 \times 8 / 48 \times 32 = 266.(6)\%$ .

**Оптимальное кол-во нитей в блоке (для Fermi):**  $48 \times 32 / 8 = 192$

## Профилирование метрик

```
.../Lab3> nvprof -m achieved_occupancy ./Lab3
```

```
==11929== NVPROF is profiling process 11929, command: ./Lab3
```

```
gTest1 took 3.56662
```

```
gTest2 took 2.41971
```

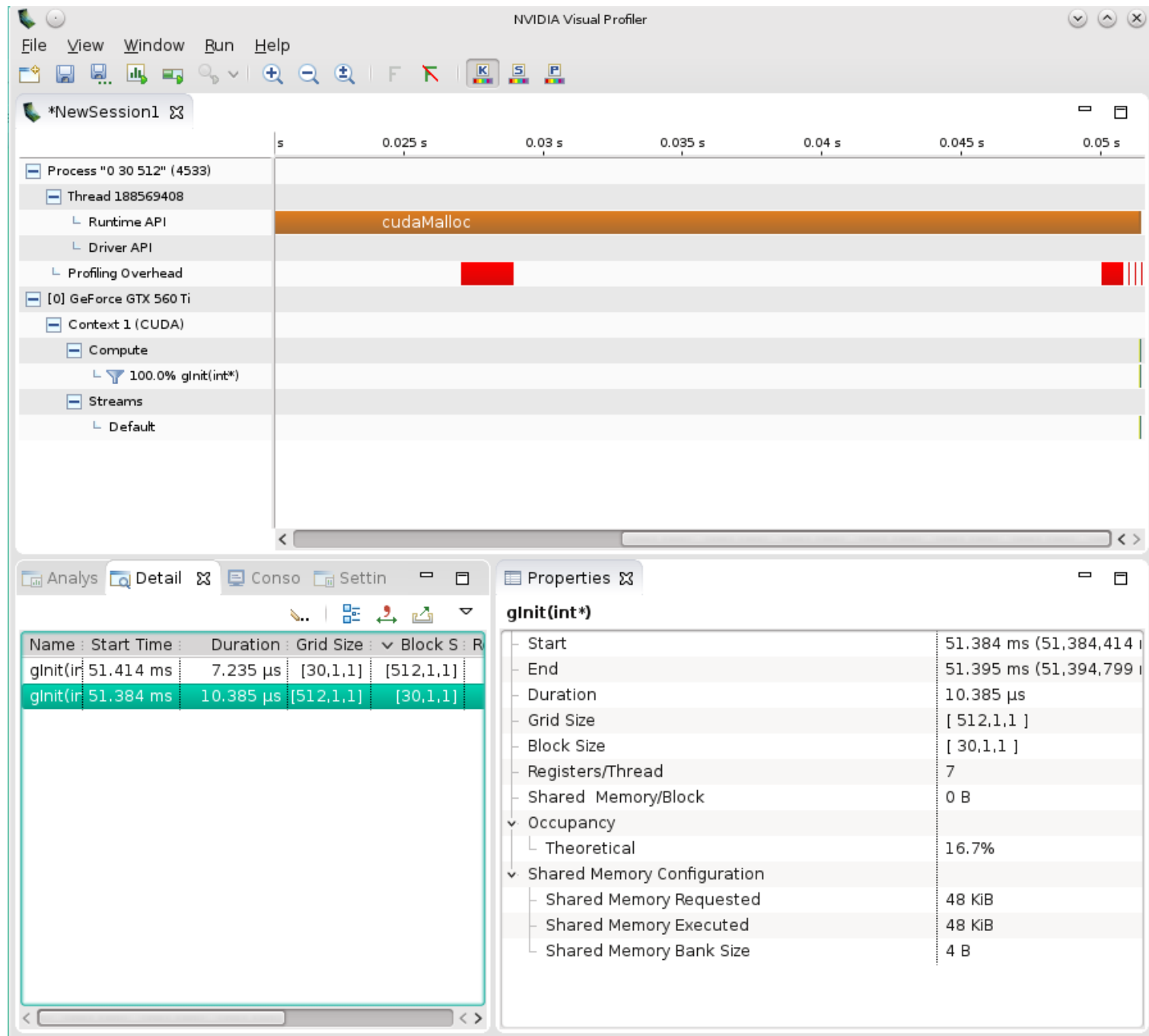
```
==11929== Profiling application: ./Lab3
```

```
==11929== Profiling result:
```

```
==11929== Metric result:
```

| Invocations                     |                    | Metric Name |          |                    | Metric Description |
|---------------------------------|--------------------|-------------|----------|--------------------|--------------------|
|                                 | Min                | Max         | Avg      |                    |                    |
| Device "GeForce GTX 560 Ti (0)" |                    |             |          |                    |                    |
| Kernel: gTest1(float*)          |                    |             |          |                    |                    |
| 1                               | achieved_occupancy |             |          | Achieved Occupancy |                    |
|                                 | 0.702714           | 0.702714    | 0.702714 |                    |                    |
| Kernel: gTest2(float*)          |                    |             |          |                    |                    |
| 1                               | achieved_occupancy |             |          | Achieved Occupancy |                    |
|                                 | 2.426397           | 2.426397    | 2.426397 |                    |                    |

>nvvp ./Lab3



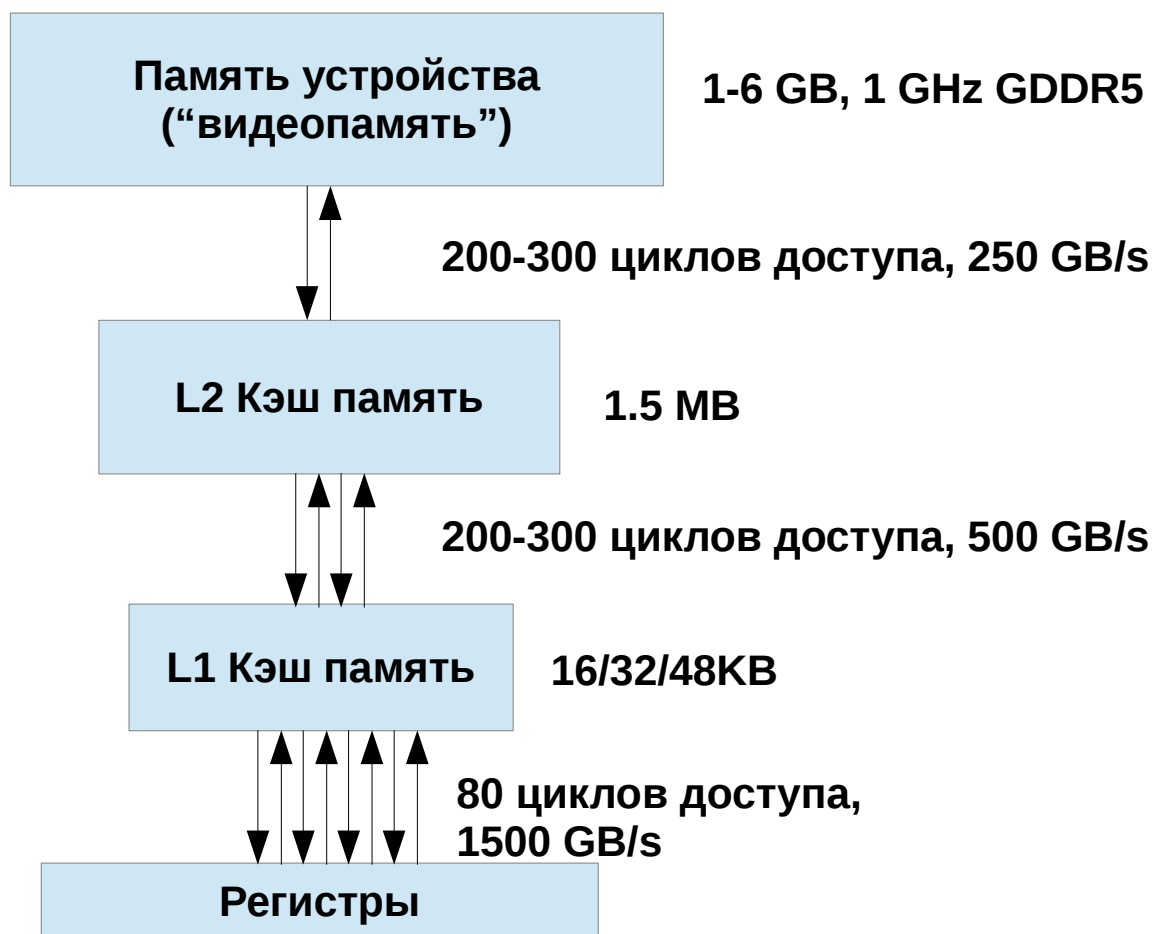


## Определение реальной заполняемости GPU

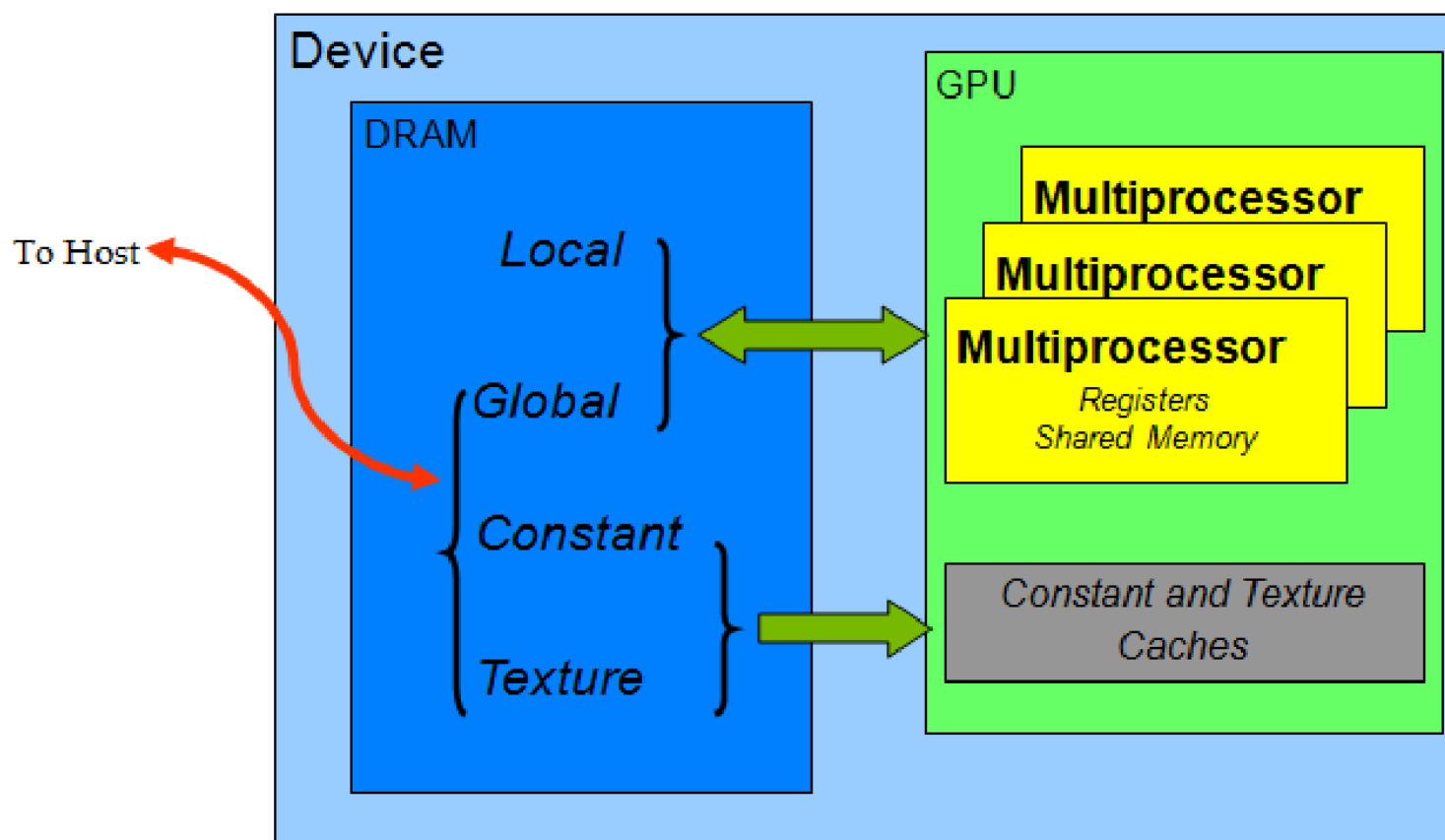
**Table 14. Technical Specifications per Compute Capability**

|   | Compute Capability |      |      |        |       |       |       |       |       |       |       |
|---|--------------------|------|------|--------|-------|-------|-------|-------|-------|-------|-------|
| Technical Specifications                              | 3.0                | 3.2  | 3.5  | 3.7    | 5.0   | 5.2   | 5.3   | 6.0   | 6.1   | 6.2   | 7.0   |
| Maximum number of resident blocks per multiprocessor  | 16                 |      |      |        | 32    |       |       |       |       |       |       |
| Maximum number of resident warps per multiprocessor   | 64                 |      |      |        |       |       |       |       |       |       |       |
| Maximum number of resident threads per multiprocessor | 2048               |      |      |        |       |       |       |       |       |       |       |
| Number of 32-bit registers per multiprocessor         | 64 K               |      |      | 128 K  | 64 K  |       |       |       |       |       |       |
| Maximum number of 32-bit registers per thread block   | 64 K               | 32 K | 64 K |        |       |       | 32 K  | 64 K  | 32 K  | 64 K  |       |
| Maximum number of 32-bit registers per thread         | 63                 | 255  |      |        |       |       |       |       |       |       |       |
| Maximum amount of shared memory per multiprocessor    | 48 KB              |      |      | 112 KB | 64 KB | 96 KB | 64 KB | 96 KB | 64 KB | 96 KB |       |
| Maximum amount of shared memory per thread block      | 48 KB              |      |      |        |       |       |       |       |       |       | 96 KB |

## Иерархия памяти графического процессора (Fermi, Kepler)

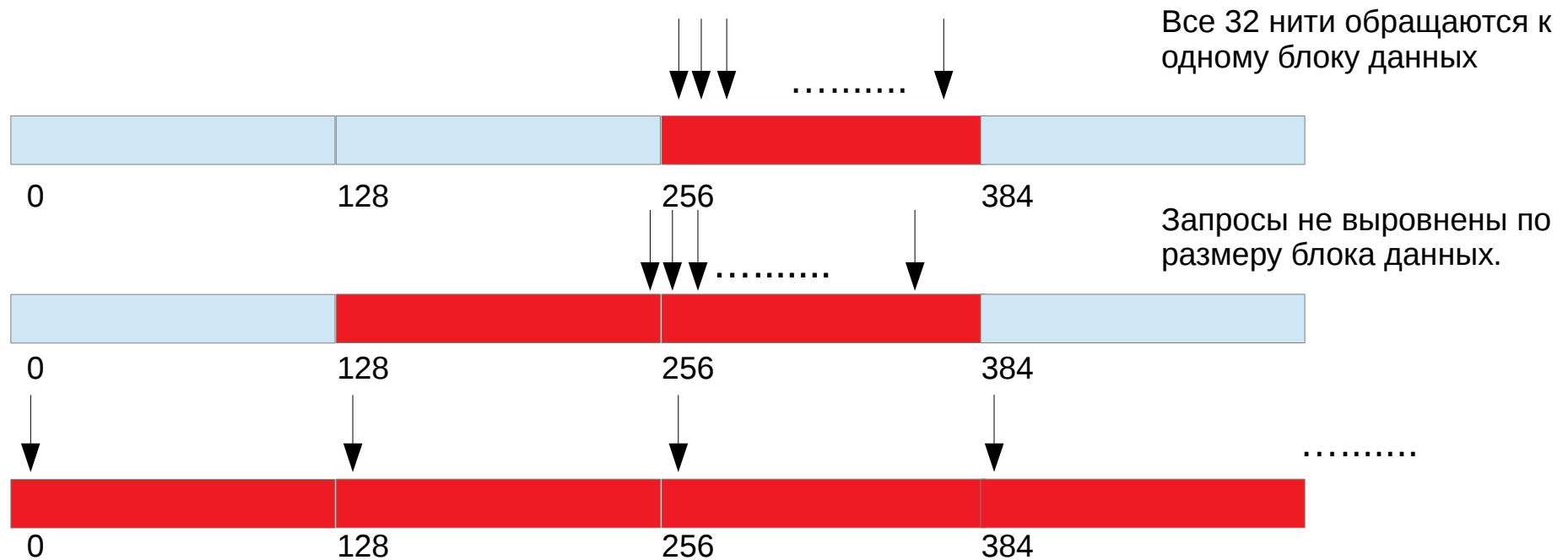


## Иерархия памяти устройства CUDA



## Объединение запросов к глобальной памяти (*coalescing*)

Запросы на чтение и запись к глобальной памяти нитями одного варпа (*a warp*) объединяются в транзакции, количество которых равно количеству необходимых для выполнения запросов блоков данных (*cache lines*) L1 кэша размером в 128 байт. **COMPUTE CAPABILITIES 2.\* !**



## *Нарушение coalescing'a (пример)*

```
__global__ void gTest1(float* a){
    int i=threadIdx.x+blockIdx.x*blockDim.x;
    int j=threadIdx.y+blockIdx.y*blockDim.y;
    int l=gridDim.x*blockDim.x;
    //int J=gridDim.y*blockDim.y;
    a[i+j*l]=(float)(threadIdx.x+blockDim.y*blockIdx.x);
}

__global__ void gTest2(float* a){
    int i=threadIdx.x+blockIdx.x*blockDim.x;
    int j=threadIdx.y+blockIdx.y*blockDim.y;
    //int l=gridDim.x*blockDim.x;
    int J=gridDim.y*blockDim.y;
    a[j+i*J]=(float)(threadIdx.y+threadIdx.x*blockDim.y);
}
```

**Спасибо за внимание**