

Оптимизация

Передача данных Host / Device

- Является дорогостоящей операцией
- Возможна асинхронная передача

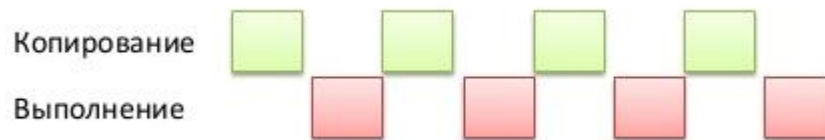
`cudaMemcpy()`

`cudaMemcpyAsync()`

- Синхронная передача



- Асинхронная передача



Типы памяти устройства

- Streaming Multiprocessor
 - Регистровая память
 - Разделяемая память
 - Константная память
- DDRAM
 - Локальная память
 - Глобальная память
- Тектурный блок
 - Тектурная память

Объединенное чтение DRAM

- Выравнивание исходных данных по границе слова
- Потоки warp-а должны осуществлять одновременное чтение DDRAM

Разделяемая память и конфликты

- Общая для всех потоков блока
- Распределяется между блоками
- Разбивается на банки (32-битные слова)

Float vs Double

- Арифметика с float числами быстрее чем с double числами

Использовать “f” при объявлении числовых констант (3.14f)

Умножение / степень

- По возможности использовать оператор сдвига
- Для известных целых значений степеней использовать явное перемножение вместо вызова `pow()`

Аппаратные аналоги функций

__sinf()	sinf()
__cosf()	cosf()
__expf()	expf()

...

Операторы ветвления

Операторы управления потоком команд (if, switch, for, while, do-while) отрицательно сказываются на производительности

следить за исполнением внутри варпа, в идеале все нити идут по одному пути, иначе возможно последовательное выполнение

Процесс оптимизации

Определить что ограничивает
производительность

- Пропускная способность памяти
- Исполнение инструкций
- Задержка
- Комбинации

Процесс оптимизации

Исследовать ограничители в порядке
важности

- Насколько эффективно работает
- Анализ и нахождение проблем
- Применить оптимизации

Анализ через модификацию кода

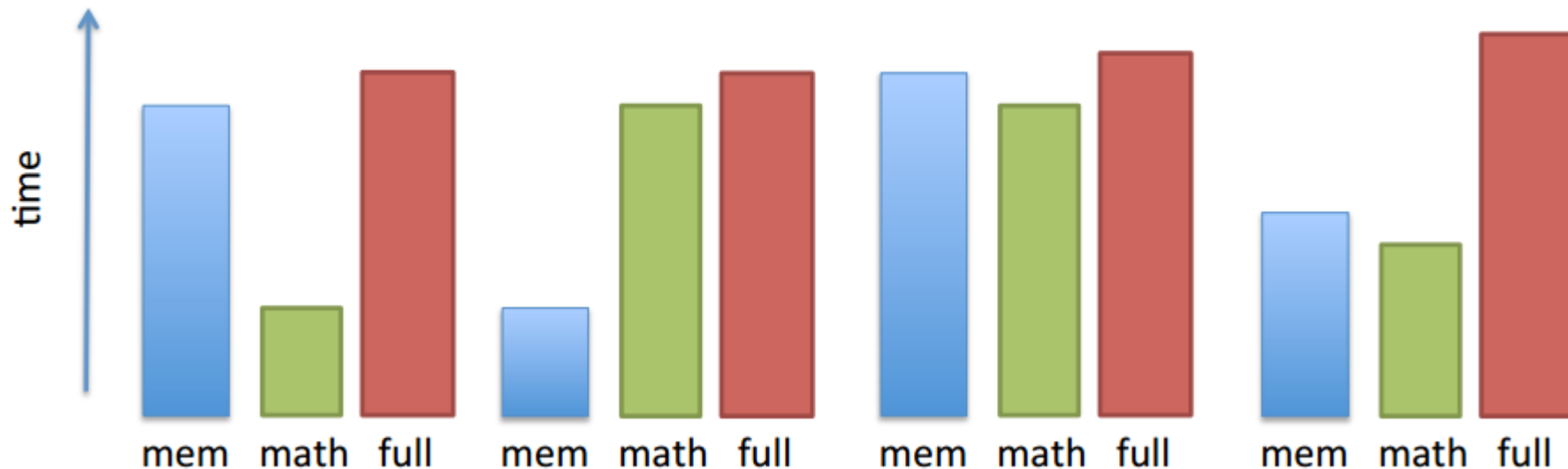
Измеряем время работы 2-х модификаций:

- Минимум инструкций, только доступ к памяти(mem)
- Минимум обращений к памяти, только инструкции(math)

Удобно для алгоритмов не использующих зависимость по данным для ветвления или адресации

Сравниваем время работы модифицированных ядер

- Помогает определить во что мы упираемся
- Насколько хорошо идет перекрытие арифметических операций и чтений из памяти



Ограничитель –
память

Хорошее
перекрытие
математики и
памяти:

нет проблем с
задержками

Ограничитель –
инструкции

Хорошее
перекрытие
математики и
памяти:

нет проблем с
задержками

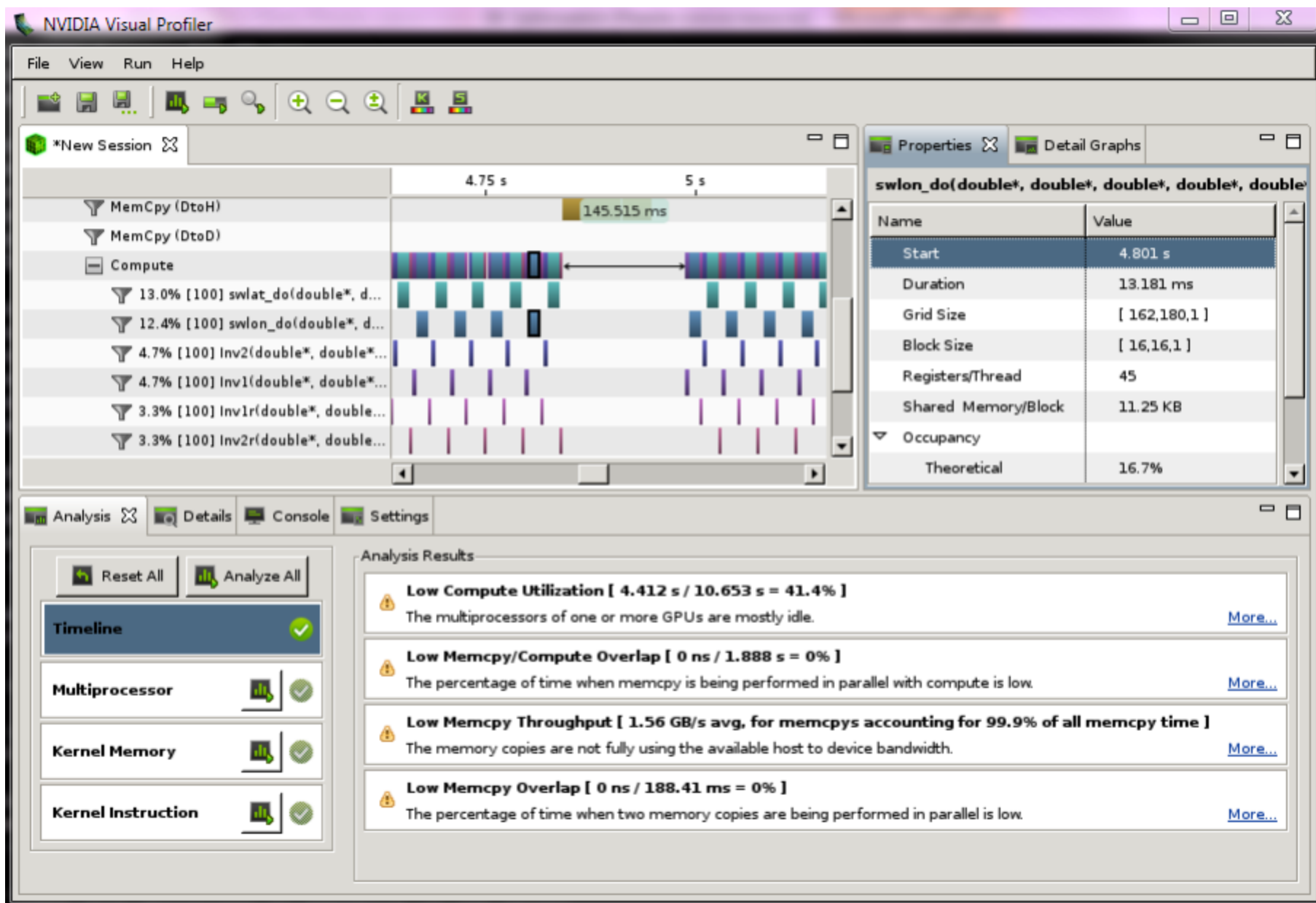
Баланс

Хорошее
перекрытие
математики и
памяти:

нет проблем с
задержками

Плохое перекрытие
математики и
памяти:

проблемы с
задержками в
время



Visual Studio interface showing CUDA code, device summary, and memory dump.

Code Snippet:

```
// Step size used to iterate through the sub-matrices of A
int aStep = BLOCK_SIZE;

// Index of the first sub-matrix of B processed by the block
int bBegin = BLOCK_SIZE * bx;

// Step size used to iterate through the sub-matrices of B
int bStep = BLOCK_SIZE * by;

// Csub is used to store the element of the block sub-matrix
// that is computed by the thread
float Csub = 0;

// Loop over all the sub-matrices of A and B
// required to compute the block sub-matrix
for (int a = aBegin, b = bBegin;
     a <= aEnd;
     a += aStep, b += bStep) {

    // Declaration of the shared memory array A
    // store the sub-matrix of A
    __shared__ float Aa[BLOCK_SIZE][BLOCK_SIZE];
```

NVIDIA Parallel Ntght - CUDA Focus Picker

Block: 4, 0, 0 Dimensions: 8, 5, 1
Thread: 14, 0, 0 16, 16, 1

Examples
#129 for block index 129
10 for coordinates 10, 0
10, 5 for coordinates 10, 5

Memory 1

Address	0x00113c00	Columns: Auto
0x00113c00	0.17999917	0.34229662
0x00113c08	0.89681691	0.36232613
0x00113c10	0.10892056	0.33661517
0x00113c18	0.71721848	0.30534621
0x00113c20	0.27188939	0.33933069
0x00113c28	0.32767726	0.047171236
0x00113c30	0.70702231	0.56959748
0x00113c38	0.32759087	0.46246061
0x00113c40	0.79228126	0.30661395
0x00113c48	0.11164815	0.048829616
0x00113c50	0.84851474	0.27295756
0x00113c58	0.34851530	0.87963499
0x00113c60	0.74705223	0.73012463
0x00113c68	0.13171728	0.85436536
0x00113c70	0.21299867	0.76390890
0x00113c78	0.60927163	0.87200534
0x00113c80	0.88998287	0.60104206
0x00113c88	0.84161735	0.18988646
0x00113c90	0.31659250	0.18057832
0x00113c98	0.98294014	0.82775354
0x00113ca0	0.72755181	0.97006525
0x00113ca8	0.28931546	0.59645373
0x00113cb0	0.30326244	0.45713779
0x00113cb8	0.0097964415	0.33207190

Locals

Name	Value	Type
blockIdx	{x = 4, y = 0, z = 0}	const int
blockDim	{x = 16, y = 16, z = 1}	const dim
gridDim	{x = 8, y = 5, z = 1}	const dim
a	???	int
b	???	int
bx	4	int
by	0	int
tx	14	int
ty	0	int
aBegin	0	int
aEnd	47	int
aStep	16	int
bBegin	64	int
bStep	2048	int
Csub	0	float
c	???	int
C	0x00113c00 0	_device_
A	0x00113000 0.20208646	_device_
B	0x00113c00 0.80645362	_device_
mA	48	_shared_
nB	128	_shared_

Счетчики профайлера

- timestamp
- gridsize
- threadblocksize
- dynsmemperblock
- stasmemperblock
- regperthread
- memtransferdir
- memtransfersize
- Streamid
- gputime
- cputime
- Occupancy
- gld_incoherent
- gld_coherent
- gld_32b / gst_32b
- gld_64b / gst_62b
- gld_128b / gst_128b
- gld_request
- gst_incoherent
- gst_coherent
- gst_request
- local_load
- local_store
- branch
- divergent_branch
- instructions
- warp_serialize
- cta_launched

- Значение имеют не цифры, а их приращение и соотношение
- Для ядер нужно стремиться чтобы стремились к нулю непоследовательные обращения к памяти

(gld_incoherent, gld_coherent, gst_incoherent, gst_coherent)

Арифметика - Память

- Оптимальное соотношение для Tesla C2050
инструкции : байты
 - ~ 3.6 : 1, float при включенном ECC
 - ~ 4.5 : 1, float при выключенном ECC

Счетчик профилировщика (на варп)

— сколько инструкций исполнилось

Instructions executed

— Конфликты банков разделяемой памяти

l1_shared_bank_conflict,

shared_load, shared_store

Конфликты существенны если:

$l1_shared_bank_conflict \gg (shared_load + shared_store)$

$l1_shared_bank_conflict \gg instructions_issued$

Спиллинг регистров

Достигнут предел доступных регистров,
компилятор начинает использовать
локальную память

- Дополнительная нагрузка шины памяти
- Дополнительные инструкции
- Не всегда проблема

счетчики l1_local_load_hit, l1_local_load_miss

CUDA Occupancy calculator

Just follow steps 1, 2, and 3 below! (or click here for help)		
1.) Select Compute Capability (click):		1,2
Physical Limits for GPU:	2.) Enter your resource usage:	
Threads / Warp	Threads Per Block	256
Warps / Multiprocessor	Registers Per Thread	18
Threads / Multiprocessor	Shared Memory Per Block (bytes)	512
Thread Blocks / Multiprocessor		
Total # of 32-bit registers / Multiprocessor		
Shared Memory / Multiprocessor (bytes)	(Don't edit anything below this line)	
Allocation Per Thread Block	3.) GPU Occupancy Data is displayed here and in the graphs:	
Warps	Active Threads per Multiprocessor	768
Registers	Active Warps per Multiprocessor	24
Shared Memory	Active Thread Blocks per Multiprocessor	3
These data are used in computing the occupancy	Occupancy of each Multiprocessor	75%
Maximum Thread Blocks Per Multiprocessor	Blocks	
Limited by Max Warps / Multiprocessor	4	
Limited by Registers / Multiprocessor	3	
Limited by Shared Memory / Multiprocessor	32	
Thread Block Limit Per Multiprocessor highlighted	RED	