# Mathematical computations with GPUs

## CUDA program optimization

Alexey A. Romanenko

arom@ccfit.nsu.ru
Novosibirsk State University

# Workflow

* Gathering information
  * memory bandwidth
  * instruction throughput
  * latencies
  * all above
* Studying limitation factors ordered by priority
  * Measure
  * Analyses
  * Optimization

# Profiling

* Environment variables (old style)
    * CUDA_PROFILE=1
    * CUDA_PROFILE_LOG="cuda_profile_%p_%d.log "
    * CUDA_PROFILE_CONFIG
    * CUDA_PROFILE_CSV=1
* Environment variables (modern style)
    * COMPUTE_PROFILE …

# Profiling. Configure file
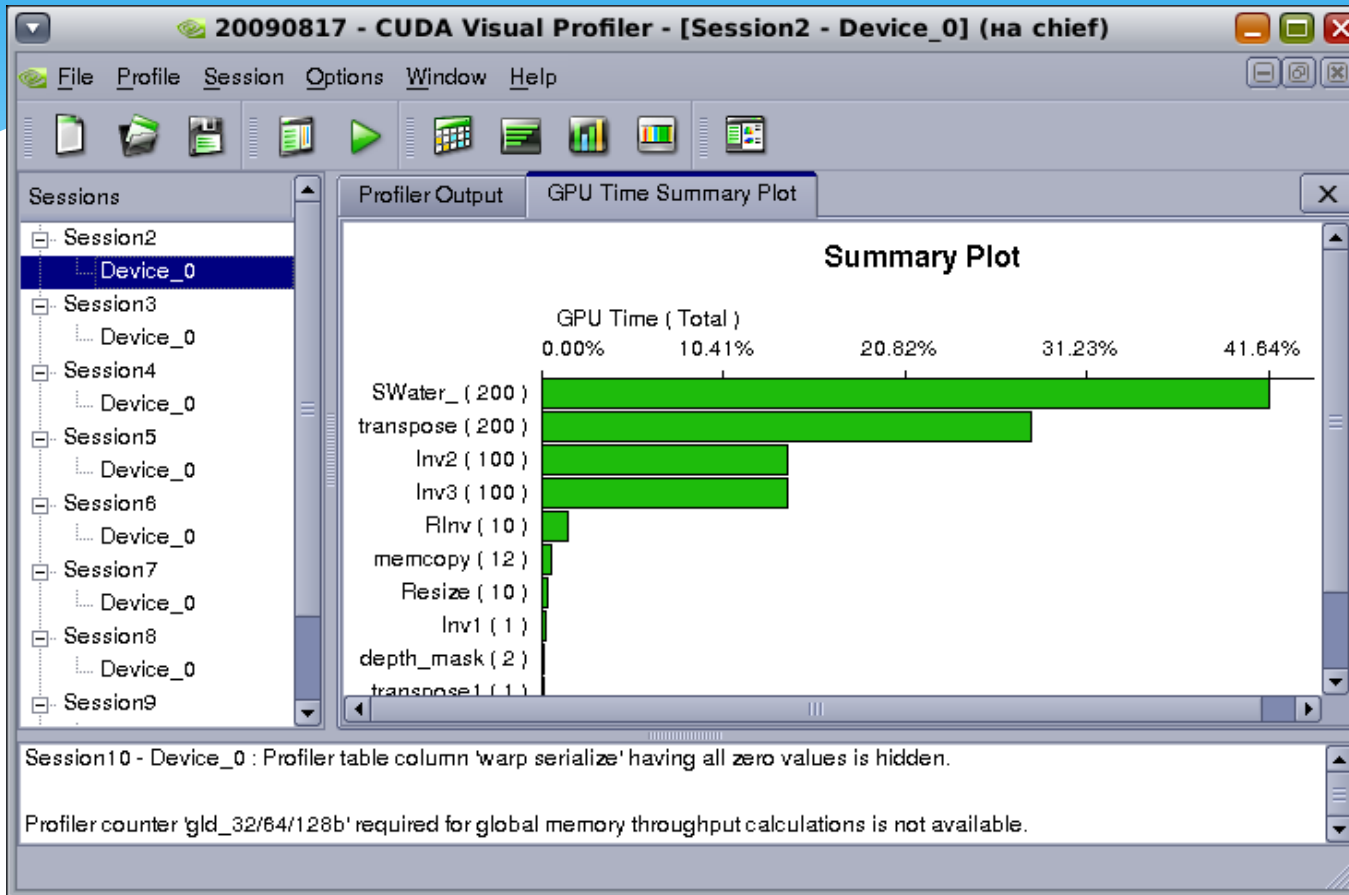
* timestamp
* gridsize
* threadblocksize
* dynsmemperblock
* stasmemperblock
* regperthread
* memtransferdir
* memtransfersize
* streamid

# Profiling.
# Configure file(CUDA 2.3)

* 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
* gputime
* cputime
* occupancy

# Profiling. Output

```
method,gputime,cputime,occupancy
method=[ memcopy ] gputime=[ 5519.680 ]
method=[ memcopy ] gputime=[ 5516.992 ]
method=[ memcopy ] gputime=[ 5517.728 ]
method=[ memcopy ] gputime=[ 5508.288 ]
method=[ __globfunc__Z10depth_maskPfS_fii ] gputime=[ 1791.424 ] cputime=[ 14.000 ]
                                              occupancy=[ 0.812 ]
method=[ __globfunc__Z10depth_maskPfS_fii ] gputime=[ 1769.920 ] cputime=[ 2.000 ]
                                              occupancy=[ 0.812 ]
method=[ __globfunc__Z9transposePfS_ii ] gputime=[ 1598.336 ] cputime=[ 2.000 ]
                                              occupancy=[ 0.812 ]
method=[ __globfunc__Z12Invariants_XPfS_S_S_S_S_ii ] gputime=[ 5061.920 ]
                      cputime=[ 3.000 ]   occupancy=[ 0.812 ]
method=[ __globfunc__Z7SWater_PfS_S_S_S_S_fS_S_S_ii ] gputime=[ 10506.752 ]
                      cputime=[ 2.000 ]   occupancy=[ 0.406 ]
method=[ __globfunc__Z9transposePfS_ii ] gputime=[ 1599.776 ] cputime=[ 2.000 ]
                                              occupancy=[ 0.812 ]
method=[ __globfunc__Z9transposePfS_ii ] gputime=[ 1613.472 ] cputime=[ 1.000 ]
                                              occupancy=[ 0.812 ]
method=[ __globfunc__Z9transposePfS_ii ] gputime=[ 1609.056 ] cputime=[ 2.000 ]
                                              occupancy=[ 0.812 ]
```

# Cudaprof

# Computerprof (CUDA 4.0)

* New user' interface

* Improvements

  * Gives recommendations based on gathered information

  * Automatically detects limitation factors

# NVVP (CUDA 4.1)

# NVIDIA Parallel Nsight

# NSight Eclipse Edition

# Profiling. CUDA+MPI

```
# Open MPI
if [ ! -z ${OMPI_COMM_WORLD_RANK} ] ; then
rank=${OMPI_COMM_WORLD_RANK}
fi
# MVAPICH
if [ ! -z ${MV2_COMM_WORLD_RANK} ] ; then
rank=${MV2_COMM_WORLD_RANK}
fi
# INTEL
if [ ! -z ${PMI_RANK} ] ; then
rank=${PMI_RANK}
fi
# Set the nvprof command and arguments.
NVPROF="nvprof --output-profile outfile.$rank $nvprof_args"
exec $NVPROF $*
```

# Profiling. CUDA+MPI

```
mpirun -np 4 ./nvprof-script.sh --print-api-trace ./a.out


CUDA 5.5
mpirun -np 2 nvprof --output-profile output.%p ./a.out


nvprof -i out.14895
======= Profiling result:
Time(%)      Time     Calls       Avg       Min       Max  Name
 50.26%  3.3920us         2  1.6960us  1.5040us  1.8880us  [CUDA memcpy HtoD]
 29.87%  2.0160us         1  2.0160us  2.0160us  2.0160us  [CUDA memcpy DtoH]
 19.87%  1.3410us         1  1.3410us  1.3410us  1.3410us  AddVectors(...
```

# Analysis of profiling results

* Absolute values tells you nothing
* Take into account ratio of counters and how value of a counter changes.
  * gld_incoherent, gst_incoherent  moved to zero.

# Instructions or memory

* Optimal value instruction:byte for Tesla C2050:
    * ~3.6 : 1, float, ECC on
    * ~4.5 : 1, float, ECC off
* counters
    * 32*instructions_issued (+1 for warp)
    * 128B*(global_store_transaction+ l1_global_load_miss) (+1 one cache line of L1)
* CUDA 4.0+ defines limiters automatically

# Instruction analysis

* performance counters (per warp)
  * instructions executed: number of executed instructions
  * instructions issued: including serelization
* Bigger difference – bigger problems. Cache misses.
* Optimal values for the device
  * see Programming Guide or Visual Profiler

# Serialization

* Warp divergence
  * Counters: divergent_branch, branch
  * Percent of divergent branches
* Shared bank conflicts
  * Conters:
    * l1_shared_bank_conflict,
    * shared_load, shared_store
* Conflicts are obstacle if both conditions are true:
  * l1_shared_bank_conflict >> (shared_load+ shared_store)
  * l1_shared_bank_conflict >> instructions_issued
* CUDA 4.0+ defines limiters automatically

# Register spilling

* Compiler could move values from registers to local memory (spilling)
    * Fermi has 63register per thread max.
    * User could define maximum number of registers
    * Local memory works as global but cached in  L1
        * L1 cache miss – request to global memory
    * Compiler flag –ptxas-options=–v shows number of local memory, number of registers for kernels
* Could influence on performance
    * Additional memory traffic
    * Additional instructions
    * Could not de a problem

# Register spilling

* counters: l1_local_load_hit, l1_local_load_miss
    * Influence on # of instructions
    * Compare with total number of instructions
* Influence on memory band width
    * Compare $2*$l1_local_load_miss with global memory operation (read + write)

# Occupancy

* Occupancy - the ratio of the number of active warps per multiprocessor to the maximum number of warps that can be active on the multiprocessor at once

# CUDA Occupancy calculator

**1.) Select Compute Capability (click):** | **1,2**

**Physical Limits for GPU:**

| | |
|---|---|
| Threads / Warp | |
| Warps / Multiprocessor | |
| Threads / Multiprocessor | |
| Thread Blocks / Multiprocessor | |
| Total # of 32-bit registers / Multiprocessor | |
| Shared Memory / Multiprocessor (bytes) | |

**2.) Enter your resource usage:**

| | |
|---|---|
| Threads Per Block | 256 |
| Registers Per Thread | 18 |
| Shared Memory Per Block (bytes) | 512 |

(Don't edit anything below this line)

**Allocation Per Thread Block**

| | |
|---|---|
| Warps | |
| Registers | |
| Shared Memory | |
| These data are used in computing the o... | |

**3.) GPU Occupancy Data is displayed here and in the graphs:**

| | |
|---|---|
| Active Threads per Multiprocessor | 768 |
| Active Warps per Multiprocessor | 24 |
| Active Thread Blocks per Multiprocessor | 3 |
| 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 |

# CUDA Occupancy calculator

# CUDA occupancy API

* Appeared in CUDA 6.5
* <CUDA_Toolkit_Path>/include/cuda_occupancy.h
    * cudaOccupancyMaxActiveBlocksPerMultiprocessor()
    * cudaOccupancyMaxPotentialBlockSize()
    * cudaOccupancyMaxPotentialBlockSizeVariableSMem()

* Example:
    * http://devblogs.nvidia.com/parallelforall/cuda-pro-tip-occupancy-api-simplifies-launch-configuration/

# Optimization

* Instruction
    * Load of operands
    * Execute instruction
    * Store result
* Optimization
    * Use 'fast' instructions
    * Reduce memory access latency
    * Fully utilize memory band width

# Instructions

* Arithmetic operations (cc = 1.2):
    * 4 clocks - FMUL, FADD, FMAD IADD, binary operations, ICMP, MIN, MAX
    * 16 clocks - __log, 1/sqrt, IMUL, 1/(float)x
    * 32 clocks - sqrt, __sin, __cos, __exp
    * 36 clocks FDIV
    * 20 clocks - __fdividef(x, y)
* Compiler options:
    * -ftz=true
    * -prec-div=false
    * -prec-sqrt=false

# If conditions. Branching

* Existing of two execution paths in a warp can cause threads to diverge (i.e. to follow different execution paths).

* Minimize flow control instructions (within a warp).
  * Pre-calculations
  * Re-arrangement of threads

# Memory access

* Global memory is accessed via 32-, 64-, or 128-byte memory transactions.
* Only the 32-, 64-, or 128-byte segments of device memory that are aligned to their size can be read/written by memory transactions.
* To maximize global memory throughput, it's important to maximize coalescing by:
    * Following the most optimal access patterns,
    * Using data types that meet the size and alignment requirement,
    * Padding data in some cases.

* Use shared, constant memory, and textures.

# Memory access

* Use cudaMalloPitch for 2D arrays
* Use CUDA arrays (cudaMallocArray) for 2D , 3D arrays
* Use page-locked memory for Device-Host operations
    * cudaHostAlloc(), cudaFreeHost(), cudaHostRegister(),
* Use textures (massy access pattern)
* Use surfaces (CUDA 4.0)

# L1 – caching and size (Fermi)

* possibilities:
  * L1 cache on
    * By default (option -Xptxas –dlcm=ca)
    * Memory transaction - 128 bytes
  * L1 cache off
    * option –Xptxas –dlcm=cg
    * Memory transaction - 32 bytes
* Cache size (L1/SMEM)
  * 16KB L1, 48KB SMEM or 48KB L1, 16KB SMEM
  * Set with CUDA API
* Recommendation:
  * Try all possibilities (CA, CG) x(16, 48)

# Concurrent kernel execution

* Available if:
    * Device with computer compatibility greater than 2.0
    * Device prop. concurrentKernels == 1
    * Kernels are from the same context
* Max number of concurrent kernels - 16

# Data transfer operations

* D2D memory bandwidth is much higher than H2D and D2H
* It could give benefits to  launch kernel with low parallelism than to copy data to Host, process them and return back to Device.
* Use page-locked memory cudaMallocHost()
    * Async operations
    * Overlap to/from data transfer operations
* Overlap data transfer operation with kernel execution

# Concurrent data transfer

* Possible if:
    * page-locked memory is used

    * Device computer compatibility >= 2.0
    * Property **asyncEngineCount** = 2

# Overlapping data transfer with kernel execution

* Possible if:
    * Device computer compatibility) >= 1.1
    * Property asyncEngineCount > 0
* Not allowed with CUDA Arrays or 2Darrays, allocated with  cudaMallocPitch()
* Variable CUDA_LAUNCH_BLOCKING set to "1" block the possibility

# Compute capatibility

* Compute Capability 1.0+
  * Async kernel execution
* Compute Capability 1.1+ ( e.g., C1060, cc1.3 )
  * One copy engine added. Property **asyncEngineCount**
* Compute Capability 2.0+ ( e.g., C2050 )
  * Possibility of concurrent kernel execution added (property **concurrentKernels**)
  * Second copy engine added. Property **asyncEngineCount**

# GPU to GPU copying

**CUDA 3.2**

cudaMemcpy(Host, GPU1);
cudaMemcpy(GPU2, Host);

**CUDA 4.0**

cudaMemcpy(GPU1, GPU2);

Requirements:

Tesla 20xx (Fermi)

64-bit application and OS

# Registers spilling

* Change max limit for registers
    * Use __launch_bounds__

```
__global__ void
__launch_bounds__(maxThreadsPerBlock, minBlocksPerMS)
MyKernel(...){...}
```

* turn L1 cache off

* Increase  L1 cache size upto 48KB

# Registers spilling

```
arom@cuda:~/cuda/edison$ /usr/local/cuda/bin/nvcc -
gencode=arch=compute_20,code=\"sm_20,compute_20\" -m32 --compiler-
options -fno-strict-aliasing  -I. -I /usr/local/cuda/include -I
/usr/local/cudasdk/C/common/inc -I/usr/local/cudasdk/shared/inc -
DUNIX -O2 -g --ptxas-options=-v,-abi=no  -c most_cuda.cu -
maxrregcount=32

ptxas info    : Compiling entry function
'_Z8swlon_doPdS_S_S_S_S_jiiddi' for 'sm_20'
ptxas info    : Used 32 registers, 52+0 bytes lmem, 11520+0 bytes
smem, 92 bytes cmem[0], 32 bytes cmem[14], 48 bytes cmem[16]
ptxas info    : Compiling entry function
'_Z4Inv1PdS_S_S_S_S_S_jiid' for 'sm_20'
ptxas info    : Used 19 registers, 80 bytes cmem[0], 32 bytes
cmem[14], 16 bytes cmem[16]
```