Makespan computation for GPU threads running on a single streaming multiprocessor



Supported by the REGAIN project, ref. FCOMP-01-0124-FEDER-020447,

co-funded by National Funds through the FCT-MCTES ERDF through COMPETE under PhD grant SFRH/BD/82069/2011

Stream Processing

Streams Collection of data. All data is expressed in streams.

Kernels Series of operations. Input: streams. Output: streams.

Why Streams?

Data parallelism Stream elements can be processed at once.

Task parallelism Pipeline.



GPU software application

Large data collections.

Data parallelism.

High arithmetic intensity.

Minimal dependency between data elements.

Application areas





© nvidia.com





GPU as a co-processor



© Kirk, David B. and Hwu, Wen-mei W.

Related work

G. Elliott and J. Anderson.
Real-time multiprocessor systems with GPUs.
18th International Conference on Real-Time and Network Systems (RTNS), 2010.

S. Kato, K. Lakshmanan, A. Kumar, Y. Ishikawa, and R. Rajkumar. **Rgem: A responsive GPGPU execution model for runtime engines.** 32nd IEEE Real-Time Systems Symposium (RTSS), 2011.

R. Mangharam and A. A. Saba.
Anytime algorithms for GPU architectures.
32nd IEEE Real-Time Systems Symposium (RTSS), 2011.

Fermi Architecture



© nvidia.com

16 streaming multiprocessors

Streaming multiprocessor



CUDA cores

load/store units

Fermi Streaming Multiprocessor (SM)

© nvidia.com

Computation entities

Thread



thread ID



© mohaircarpets.co.za

Warp



a group of threads with consecutive IDs; all should be executed in parallel.

Control flow divergence

global__ void kv(int* x, int* y) int i = threadIdx.x + blockDim.x * blockIdx.x; int t; **bool** b = f(x[i]);if (b) // g(x) t = g(x[i]); else // h(x) t = h(x[i]));y[i] = t;}



© nvidia.com

How does it work?



float md = $(x - x_1)^2 + (y - y_1)^2$; int mdp = 1; for (int i=2; i<=N; i++) if ($(x - x_i)^2 + (y - y_i)^2 < md$){ md = $(x - x_i)^2 + (y - y_i)^2$; mdp = i; }

//minimal distance
//minimal distance point

Parallel Thread Execution (PTX)



CUDA cores

load/store units

PTX Assembly Code

	mov.u32	\$r0,	N ado	lr;		// N
	mov.f32	\$f1,	x add	r,		// x
	mov.f32	\$f2,	y_add	r;		// y
	mov.f32	\$f3,	x1_ad	dr		// x ₁
	mov.f32	\$f4,	y1_ad	dr;		// y ₁
	sub.f32	\$f5,	\$f1 ,	\$f3;		$// (x - x_1)$
	mul.f32	\$f6,	\$f 5,	\$f 5;		$// (x - x_1)^2$
	sub.f32	\$f 7,	\$f2,	\$f 4;		$// (y - y_1)$
	fma.f32	\$f8,	\$f 7,	\$f 7,	\$f 6;	// $(x - x_1)^2 + (y - y_1)^2$
	mov.f32	\$f 0,	\$f 8;			// md = $(x - x_1)^2 + (y - y_1)^2$
	mov.u32	\$r1,	1;			// int mdp = 1;
	mov.u32	\$r2,	2;			// for (int i=2; i<=N; i++)
Loop:	setp.gt.u32	p,	\$r2,	\$r0;		// for (int i=2; i<=N; i++)
@p	bra Done;					// for (int i=2; i<=N; i++)
	mov.f32	\$f3,	xi_ado	dr;		// x _i
	mov.f32	\$f4,	yi_ado	dr;		// y _i
	sub.f32	\$f5,	\$f1 ,	\$f3;		$// (x - x_i)$
	mul.f32	\$f6,	\$f 5,	\$f5;		$// (x - x_i)^2$
	sub.f32	\$f 7,	\$f2,	\$f 4;		// $(y - y_i)$
	fma.f32	\$f8,	\$f 7,	\$f 7,	\$f 6;	// $(x - x_i)^2 + (y - y_i)^2$
	setp.ge.f32	q,	\$f 8,	\$f 0;		// if $((x-x_i)^2+(y-y_i)^2) < md$
@q	bra If;					
	mov.f32	\$f0,	\$f 8;			// md = $(x - x_i)^2 + (y - y_i)^2$;
	mov.u32	\$r1,	\$r2;			// mdp = i;
If:						// if $((x-x_i)^2+(y-y_i)^2) < md$
	add.u32	\$r2,	\$r2,	1;		// for (int i=2; i<=N; i++)
	bra Loop;					
Done:	-					

L L L L L C С С C C C C C C C L C C C L C C С С С C C С C C

Kernel instruction string



float md = $(x - x_1)^2 + (y - y_1)^2$; int mdp = 1; for (int i=2; i<=N; i++) if ($(x - x_i)^2 + (y - y_i)^2 < md$){ md = $(x - x_i)^2 + (y - y_i)^2$; mdp = i; }

//minimal distance
//minimal distance point

LLLLCCCCCCCCCLLCCCCCCC

	Instruction Cache										
	1	Warp Schee	duler	Warp Scheduler							
		Dispatch I	Unit	Dispatch Unit							
Ē		+		+							
			Register File	(32,768 x 3	2-bit)						
Ē			*			+					
	Core	Core	Core	Core	LD/ST						
	Core	Core	Core	Core	LD/ST	SFU					
	Core	Core	Core	Core	LD/ST						
					LD/ST	SFU					
	Core	Core	Core Core	LD/ST							
ĺ	Core	Core	Core	Core	LD/ST	SEU					
	Core	Core	Core	Core	LD/ST	310					
	Core	Core	Core	Core	LD/ST						
	Core	Core	Core	Core	LD/ST LD/ST	SFU					
	Interconnect Network										
	64 KB Shared Memory / L1 Cache										
Ē	Uniform Cache										

CUDA cores	(C)
load/store units	(L)

Related work

An Adaptive Performance Modeling Tool for GPU Architectures

Sara S. Baghsorkhi, Matthieu Delahaye, Sanjay J. Patel, Willian D. Gropp, Wen-mei W. Hwu

15th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP), 2010.

A Memory-level and Thread-level Parallelism Aware GPU Architecture Performance Analytical Model,

Sunpyo Hong and Hyesoon Kim 36th International Symposium on Computer Architecture (ISCA-36), 2009.

Optimization Principles and Application Performance Evaluation of a Multithreaded GPU Using CUDA

Shane Ryoo, Christopher I. Rodrigues, Sara S. Baghsorkhi, Sam S. Stone, David B. Kirk, and Wen-mei W. Hwu 13th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, 2008.

The makespan





Worst – case Makespan Clock Cycle 3 6 8 102 4 5 7 9 Warp 1 С L L С L L Warp 2 L С L С L L

Assumptions

Every instruction takes single clock cycle.

Cache

A warp on a single SM

32-bit float instruction

Cache hit 1 clock cycle

Cache miss 400-600 clock cycles

Assumptions

Every instruction takes single clock cycle. No cache misses.

Scheduling Policy

CLLCLL 2 warps



Non-work-conserving schedule

Clock Cycle	1	2	3	4	5	6	7	8	9
Warp 1	С	L		L	С	L		L	
Warp 2		С	L		L	С	L		L

Work-conserving schedule

Assumptions

Every instruction takes single clock cycle. No cache misses. Work-conserving scheduling.

The Pessimistic Method

		Clock Cycle	1	2	3	4	5	6	7	8	9			
	CLLCLL	Warp 1	С	L		L	С	L		L				
		Warp 2		C	L		L	C	L		L			
		Clock Cycle	1	2	3	4	5	6	7	8	9	10	11	
		Warp 1	L		L		L		L	C		C]
		Warp 2		L		L		L		L	С		С	1
				1	1		1		1		1			_
		Clock Cycle	1	2	3	4	5	6	7	8	9	10	11	12
М	$-M^L \perp M^C$	Warp 1	L		L		L		L		С		С	
111	-m + m	Warp 2		L		L		L		L		С		С
				-	-	-							I	

 M^C a makespan for C-substring M^L a makespan for L-substring **Optimization Problem**

Decision variables

Usage of hardware units at particular clock cycle The pessimistic makespan – an upper bound Objective function Maximize the makespan Constraints

Capacity constraints Precedence constraints Work-conserving constraints

Integer Linear Problem (ILP)

Optimization Problem

maximize $\sum_{t=1}^{T} (t \times (LS_{W,I,t} + CC_{W,I,t}))$ subject to						
	t=1					
iterated variables	expression for constraint	number of constraints				
$\forall t$	$\sum_{w=1}^{W} \sum_{i=1}^{I} LS_{w,i,t} \le \sigma_L$	Т				
$\forall t$	$\sum_{w=1}^{W} \sum_{i=1}^{I} CC_{w,i,t} \le \sigma_C$	Т				
$\forall w = 1(W - 1)$	$\sum_{t=1}^{T} (t \times (LS_{w,I,t} + CC_{w,I,t})) \le \sum_{t=1}^{T} (t \times (LS_{W,I,t} + CC_{W,I,t}))$	W-1				
$\forall w, t$	$\sum_{i=1}^{I} LS_{w,i,t} \le 1$	$W \times T$				
$\forall w, t$	$\sum_{i=1}^{I} CC_{w,i,t} \le 1$	$W \times T$				
$\forall w, i$	$\sum_{t=1}^{T} LS_{w,i,t} = IL_i$	W imes I				
$\forall w, i$	$\sum_{t=1}^{T} CC_{w,i,t} = IC_i$	$W \times I$				
$\forall w, i = 1(I-1)$	$\sum_{t=1}^{T} (t \times (LS_{w,i,t} + CC_{w,i,t})) < \sum_{t=1}^{T} (t \times (LS_{w,i+1,t} + CC_{w,i+1,t}))$	W imes (I-1)				
$\forall t$	$E_t \ge 1 - \sigma_L + \sum_{w'=1}^W \sum_{i=1}^I LS_{w',i,t}$	Т				
$\forall t$	$E_t \times \sigma_L \le \sum_{w'=1}^W \sum_{i=1}^I LS_{w',i,t}$	Т				
$\forall w, i, t$	$\frac{1}{2}(TL_{w,i,t} + E_t) \le TL_{w,i,t} \lor E_t \le TL_{w,i,t} + E_t$	$W \times I \times T$				
$\forall t$	$G_t \ge 1 - \sigma_C + \sum_{w'=1}^W \sum_{i=1}^I CC_{w',i,t}$	Т				
$\forall t$	$G_t \times \sigma_C \leq \sum_{w'=1}^W \sum_{i=1}^I CC_{w',i,t}$	Т				
$\forall w, i, t$	$\frac{1}{2}(TC_{w,i,t}+G_t) \le TC_{w,i,t} \lor G_t \le TC_{w,i,t} + G_t$	$W \times I \times T$				

The Application



Resolving the issue of tractability



Makespan (clock cycles)

Future Work

Relaxation of the "no cache miss assumption".

Targeting the "whole" GPU.

Future Work

Fermi
48 warps
$$W = 420$$
,



Questions?

Thank You!