# GPU architecture part 2: SIMT control flow management

#### Caroline Collange Inria Rennes – Bretagne Atlantique

caroline.collange@inria.fr https://team.inria.fr/pacap/members/collange/

> Master 2 SIF ADA - 2019

# Outline

- Running SPMD software on SIMD hardware
  - Context: software and hardware
  - The control flow divergence problem
- Stack-based control flow tracking
  - Stacks
  - Counters
- Path-based control flow tracking
  - The idea: use PCs
  - Implementation: path list
  - Applications
- Software approaches
  - Use cases and principle
  - Scalarization
- Research directions

```
scale:
                                                                             test esi, esi
                                                                                  .L4
                                                                             ie
                                                                             sub
                                                                                  esi. 1
                 void scale(float a, float * X, int n)
Software
                                                                             xor
                                                                                  eax, eax
                 {
                                                                                  rdx, [4+rsi*4]
                                                                             lea
                      for(int i = 0; i != n; ++i)
                                                                        .L3:
                                                                             movss xmm1, DWORD PTR [rdi+rax]
                           X[i] = a * X[i];
                                                                             mulss xmm1, xmm0
                 }
                                                                             movss DWORD PTR [rdi+rax], xmm1
                                                                             add
                                                                                  rax, 4
                                                                                  rax, rdx
                                                                             cmp
                                                                                  .L3
                                                                             ine
                                                                        .L4:
                                                                             rep
         Architecture: sequential programming model
                                                                             ret
```



#### Hardware datapaths: dataflow execution model







6

# GPU microarchitecture: where it fits



# Programmer's view

- Programming model
  - SPMD: Single program, multiple data
  - One kernel code, many threads
  - Unspecified execution order between explicit synchronization barriers
- Languages
  - Graphics shaders : HLSL, Cg, GLSL
  - GPGPU : C for CUDA, OpenCL



Kernel

# Types of control flow

- Structured control flow: single-entry, single exit properly nested
  - Conditionals: if-then-else
  - Single-entry, single-exit loops: while, do-while, for...
  - Function call-return...

#### Unstructured control flow

- break, continue
- && || short-circuit evaluation
- Exceptions
- Coroutines
- 🔶 goto, comefrom
- Code that is hard to indent!

# Warps

 Threads are grouped into warps of fixed size





#### Control flow: uniform or divergent

 Control is *uniform* when all threads in the warp follow the same path

 Control is *divergent* when different threads follow different paths

```
Warp
x = 0;
// Uniform condition
                       T0
                           T1
                               T2
                                   T3
if(a[x] > 17) {
                        1 1
                                1
                                    1
       x = 1;
}
// Divergent condition
if(tid < 2) {
                         1
                           1
                                0
                                    \mathbf{0}
       x = 2;
}
```

Outcome per thread:

# Computer architect view

- SIMD execution inside a warp
  - One SIMD lane per thread
  - All SIMD lanes see the same instruction

- Control-flow differentiation using execution mask
  - All instructions controlled with 1 bit per lane
    - $1 \rightarrow perform instruction$
    - $0 \rightarrow do nothing$

# Running independent threads in SIMD

- How to keep threads synchronized?
  - Challenge: divergent control flow
- Rules of the game
  - One thread per SIMD lane
  - Same instruction on all lanes
  - Lanes can be individually disabled with execution mask
- Which instruction?
- How to compute execution mask?



```
x = 0;
// Uniform condition
if(tid > 17) {
      x = 1;
}
// Divergent conditions
if(tid < 2) {
      if(tid == 0) {
          x = 2;
      else {
          x = 3;
      }
}
```

Execution mask inside if statement = if condition

$\mathbf{x} = 0$ :		Warp				
<pre>// Uniform condition if(a[x] &gt; 17) {</pre>	if condition: a[x] > 17?	T0 1	T1 1	T2 1	T3 1	
x = 1;	Execution mask:	1	1	1	1	
}						
<pre>// Divergent condition if(tid &lt; 2) {</pre>	if condition: tid < 2?	1	1	0	0	
x = 2; }	Execution mask:	1	1	0	0	

# State of the art in 1804

The Jacquard loom is a GPU



- Multiple warp threads with per-thread conditional execution
  - Execution mask given by punched cards
  - Supports 600 to 800 parallel threads

# Conditionals that no thread executes

- Do not waste energy fetching instructions not executed
  - Skip instructions when execution mask is all-zeroes



Uniform branches are just usual scalar branches

# What about loops?

- Keep looping until all threads exit
  - Mask out threads that have exited the loop



### What about nested control flow?

```
x = 0;
// Uniform condition
if(tid > 17) {
      x = 1;
}
// Divergent conditions
if(tid < 2) {
      if(tid == 0) {
          x = 2;
      }
      else {
          x = 3;
      }
}
```

• We need a generic solution!

# Outline

- Running SPMD software on SIMD hardware
  - Context: software and hardware
  - The control flow divergence problem
- Stack-based control flow tracking
  - Stacks
  - Counters
- Path-based control flow tracking
  - The idea: use PCs
  - Implementation: path list
  - Applications
- Software approaches
  - Use cases and principle
  - Scalarization
- Research directions

• What do Star Wars and GPUs have in common?

# Answer: Pixar!

- In the early 1980's, the Computer Division of Lucasfilm was designing custom hardware for computer graphics
- Acquired by Steve Jobs in 1986 and became *Pixar*
- Their core product: the *Pixar Image Computer*



This early GPU handles nested divergent control flow!

#### Pixar Image Computer: architecture overview





#### The mask stack of the Pixar Image Computer



A. Levinthal and T. Porter. Chap - a SIMD graphics processor. SIGGRAPH'84, 1984.

# Observation: stack content is a histogram



- On structured control flow: columns of 1s
  - A thread active at level n is active at all levels i<n</p>
  - Conversely: no "zombie" thread gets revived at level i>n if inactive at n

# Observation: stack content is a histogram



- On structured control flow: columns of 1s
  - A thread active at level n is active at all levels i<n</p>
  - Conversely: no "zombie" thread gets revived at level i>n if inactive at n
- The height of each column of 1s is enough
  - Alternative implementation: maintain an activity counter for each thread

R. Keryell and N. Paris. Activity counter : New optimization for the dynamic scheduling of 28 SIMD control flow. ICPP '93, 1993.

# With activity counters

Code x = 0;// Uniform condition if(tid > 17) { skip x = 1;// Divergent conditions **if**(tid < 2) { inc **if**(tid == 0) { inc x = 2;} dec else { inc x = 3;} dec } dec

Counters 1 (in)activity counter / thread



# Activity counters in use

• In an SIMD prototype from École des Mines de Paris in 1993



In Intel integrated graphics since 2004



# Outline

- Running SPMD software on SIMD hardware
  - Context: software and hardware
  - The control flow divergence problem
- Stack-based control flow tracking
  - Stacks, counters
  - A compiler perspective
- Path-based control flow tracking
  - The idea: use PCs
  - Implementation: path list
  - Applications
- Software approaches
  - Use cases and principle
  - Scalarization
- Research directions

## Back to ancient history

Microsoft DirectX

7.x	8.0	8.1	9.0 a	9.0b	9.0c	10.0	10.1	11
Linified shaders								

**NVIDIA** 

	NV10	NV20	NV30	NV40		G70	G8	80-G90	GT200	) GF100
FP	16	Programmat shaders	ble FP 32	Dynamic control flo	) VV	SIMT	(	CUDA		
ATI	/AMD		FP 24			СТМ	F	=P 64     <b>(</b>	CAL	
	R100	R200	R300	R400	)	R500	)	R600	R700	Evergreen



# Early days of programmable shaders

#### It is 21<sup>st</sup> century!

• Graphics cards now look and sound like hair dryers



- Graphics shaders are programmed in assembly-like language
  - Direct3D shader assembly, OpenGL ARB Vertex/Fragment Program...
  - Control-flow: if, else, endif, while, break... are assembly instructions
- Graphics driver performs a straightforward translation to GPU-specific machine language

# Goto considered harmful?

MIPS	NVIDIA Tesla (2007)	NVIDIA Fermi (2010)	Intel GMA Gen4 (2006)	Intel GMA SB (2011)	AMD R500 (2005)	AMD R600 (2007)	AMD Cayman (2011)
j jal jr syscall	bar bra brk brkpt cal cont kil pbk pret ret ssy trap .s	bar bpt bra brk brx cal cont exit jcal jmx kil pbk pret ret ssy .s	jmpi if iff else endif do while break cont halt msave mrest push pop	jmpi if else endif case while break cont halt call return fork	jump loop endloop rep endrep breakloop breakrep continue	<pre>push push_else pop loop_start loop_start_no_al loop_start_dx10 loop_end loop_continue loop_break jump else call call_fs return return_fs alu alu_push_before alu pop after</pre>	push_else pop push_wqm pop_wqm else_wqm jump_any reactivate reactivate_wqm loop_start loop_start_no_al loop_start_dx10 loop_end loop_end loop_break jump else call call fs

Control instructions in some CPU and GPU instruction sets

alu\_pop2\_after alu\_continue alu\_break alu\_else\_after GPUs: instruction set expresses control flow structure

alu\_pop2\_after

alu\_else after

alu continue

alubreak

return

alu

return fs

alu\_push\_before
alu\_pop\_after

Where should we stop?

0

# Next: compilers for GPU code

- High-level shader languages
  - C-like: HLSL, GLSL, Cg
  - Then visual languages (UDK)
- General-purpose languages
  - CUDA, OpenCL
  - Then directive-based: OpenACC, OpenMP 4
  - Python (Numba)...



- Incorporate function calls, switch-case, && and ||...
- Demands a compiler infrastructure
  - A Just-In-Time compiler in graphics drivers

- First: turns all structured control flow into gotos, generates intermediate representation (IR)
  - e.g. Nvidia PTX, Khronos SPIR, Ilvm IR
- Then: performs various compiler optimizations on IR
- Finally: reconstructs structured control flow back from gotos to emit machine code
  - Not necessarily the same as the original source!

## Issues of stack-based implementations

If GPU threads are actual threads, they can synchronize?

- e.g. using semaphores, mutexes, condition variables...
- Problem: SIMT-induced livelock

```
while(!acquire(lock)) {}
...
release(lock)
```

Example: critical section Thread 0 acquires the lock, keeps looping with other threads of the warp waiting for the lock. Infinite loop, lock never released.

Stack-based SIMT divergence control can cause starvation!

# Issues of stack-based implementations

• Are all control flow optimizations valid in SIMT?



- What about context switches?
  - e.g. migrate one single thread of a warp
  - Challenging to do with a stack
- Truly general-purpose computing demands more flexible techniques

# Outline

- Running SPMD software on SIMD hardware
  - Context: software and hardware
  - The control flow divergence problem
- Stack-based control flow tracking
  - Stacks, counters
  - A compiler perspective
- Path-based control flow tracking
  - The idea: use PCs
  - Implementation: path list
  - Applications
- Software approaches
  - Use cases and principle
  - Scalarization
- Research directions
#### With 1 PC / thread



#### Mask stacks vs. per-thread PCs

- Before: stack, counters
  - O(n), O(log n) memory
     n = nesting depth
  - 1 R/W port to memory
  - Exceptions: stack overflow, underflow
- Vector semantics
  - Structured control flow only
  - Specific instruction sets

- After: multiple PCs
  - O(1) memory
  - No shared state
  - Allows thread suspension, restart, migration
- Multi-thread semantics
  - Traditional languages, compilers
  - Traditional instruction sets
- Can be mixed with MIMD
- Straightforward implementation is more expensive

#### Path-based control flow tracking

• A **path** is characterized by a PC and execution mask



• The mask encodes the **set of threads** that have this PC

$$\{T_{1}, T_{3}, T_{4}, T_{7}\}$$
 have PC 17

#### A list of paths represents a vector of PCs



- Worst case: 1 path per thread
  - Path list size is bounded
- PC vector and path list are equivalent
  - You can switch freely between MIMD thinking and SIMD thinking!

Select an active path









## Divergent branch is path insertion



#### Convergence is path fusion

- When two paths have the same PC, we can merge them
  - New set of threads is the union of former sets
  - New execution mask is bitwise OR of former masks



## Path scheduling is graph traversal

Degrees of freedom Active . 001001 4 0 Which path is the active path? 12 100000000 Insert At which place are new paths inserted? 01011001 17 When and where do we check for 17 000001  $\mathbf{0}$ convergence?

Compare

Different answers yield different policies

## Depth-first graph traversal

- Remember graph algorithm theory
  - Depth-first graph traversal using a stack worklist
- Path list as a stack
   = depth-first traversal of the control-flow graph
  - Most deeply nested levels first



x = 0;// Uniform condition if(tid > 17) { x = 1;} // Divergent conditions if(tid < 2) { if(tid == 0) { x = 2;} else { x = 3;} }

Question: is this the same as Pixar-style mask stack? Why?

## Breadth-first graph traversal

Goal: guarantee forward progress to avoid SIMT-induced livelocks



• Path list as a queue: follow paths in round-robin



Drawback: may delay convergence

A. ElTantawy and T. Aamodt. MIMD Synchronization on SIMT Architectures. Micro 2016

#### Limitations of static scheduling orders

- Stack works well for structured control flow
  - Convergence happens in reverse order of divergence
- But not so much for unstructured control flow
  - Divergence and convergence order do not match

#### Priority-based graph traversal

• Sort the path list based on its contents





3	0	0	1	0	0	1	1	0
12	1	0	0	0	0	0	0	0
17	0	1	0	1	1	0	0	1

Ordering paths by priority enables a scheduling policy

## Scheduling policy: min(SP:PC)

Which PC to choose as master PC ?	Source	Assembly	Order
<ul> <li>Conditionals, loops</li> <li>Order of code addresses</li> <li>min(PC)</li> <li>Functions</li> <li>Favor max nesting depth</li> </ul>	if(…) { } else { }	<pre>m m p? br else m br endif else: m endif:</pre>	
<ul> <li>min(SP)</li> <li>With compiler support</li> </ul>	while(…) { }	start:  p? br start 	
<ul> <li>Unstructured control flow too</li> <li>No code duplication</li> <li>Full backward and forward compatibility</li> </ul>	 f(); void f() {  }	 call f  f:  ret	<ul> <li>↓ ①</li> <li>↓ ③</li> <li>↓ ②</li> <li>↓ 20</li> <li>57</li> </ul>

#### Convergence with min(PC)-based policies

Sorting by PC groups paths of equal PC together

4	0	0	1	0	0	0	1	0
12	1	0	0	0	0	0	0	0
17	0	1	0	1	1	0	0	1
17	0	0	0	0	0	1	0	0

- Priority order: active path is top entry
- Convergence detection: only needed between top entry and following entries
  - No need for associative lookup

#### Example: Nvidia Volta (2017)

# Supports independent thread scheduling inside a warp

- Threads can synchronize with each other inside a warp
- Diverged threads can run barriers (as long as all threads eventually reach a barrier)



#### Advertisement: Simty, a SIMT CPU

- Proof of concept for priority-based SIMT
  - Written in synthesizable VHDL
  - Runs the RISC-V instruction set (RV32I)
  - Fully parametrizable warp size, warp count
  - 10-stage pipeline



#### Outline

- Running SPMD software on SIMD hardware
  - Context: software and hardware
  - The control flow divergence problem
- Stack-based control flow tracking
  - Stacks, counters
  - A compiler perspective
- Path-based control flow tracking
  - The idea: use PCs
  - Implementation: path list
  - Applications
- Software approaches
  - Use cases and principle
  - Scalarization
- Research directions

#### SIMT vs. multi-core + explicit SIMD

- SIMT
  - All parallelism expressed using threads
  - Warp size implementationdefined
  - Dynamic vectorization



Example: Nvidia GPUs

- Multi-core + explicit SIMD
  - Combination of threads, vectors
  - Vector length fixed at compiletime
  - Static vectorization



Example: most CPUs, Intel Xeon Phi, AMD GCN GPUs

Are these models equivalent?

## Bridging the gap between SPMD and SIMD





#### SPMD to SIMD: hardware or software ?





SIMD CPU, GPU, Xeon Phi...

- → Which is best? : open question
- $\rightarrow$  Combine both approaches?

#### Tracking control flow in software

- Use cases
  - Compiling shaders and OpenCL for AMD GCN GPUs
  - Compiling OpenCL for Xeon Phi
  - ispc: Intel SPMD Program Compiler, targets various SIMD instruction sets
- Compiler generates code to compute execution masks and branch directions
  - Same techniques as hardware-based SIMT
  - But different set of possible optimization

#### Compiling SPMD to predicated SIMD

```
x = 0;
// Uniform condition
if(tid > 17) {
      x = 1;
}
// Divergent conditions
if(tid < 2) {
      if(tid == 0) {
          x = 2;
      }
      else {
          x = 3;
      }
}
```

#### Compiling SPMD to predicated SIMD

x = 0;// Uniform condition if(tid > 17) { x = 1;} // Divergent conditions if(tid < 2) { if(tid == 0) { x = 2;} else { x = 3;}

}

(m0) mov  $x \leftarrow 0$  // m0 is current mask (m0) cmp c←tid>17 // vector comparison and m1←m0&c // compute if mask jcc(m1=0) endif1 // skip if null (m1) mov x←1 endif1: (m0) cmp c←tid<2 and m2←m0&c jcc(m2=0) endif2 (m2) cmp c←tid==0 and m3←m2&c jcc(m3=0) else (m3) mov x←2 else: and m4←m2&~c jcc(m4=0) endif2 (m4) mov x←3 endif2:

#### Benefits and shortcomings of s/w SIMT

#### Benefits

- No stack structure to maintain
  - Use mask registers directly
  - Register allocation takes care of reuse and spills to memory
- Compiler knowing precise execution order enables more optimizations
  - Turn masking into "zeroing": critical for out-of-order architectures
  - Scalarization: demoting uniform vectors into scalars

#### Shortcomings

- Every branch is divergent unless proven otherwise
  - Need to allocate mask register either way
- Restricts freedom of microarchitecture for runtime optimization

#### Scalars in SPMD code

 Some values and operations are inherently scalar

> Loop counters, addresses of consecutive accesses...

SPMD code





#### Uniform and affine vectors

- Uniform vector
  - In a warp, v[i] = c
  - Value does not depend on lane ID

- Affine vector
  - In a warp, v[i] = b + i s
  - Base b, stride s
  - Affine relation between value and lane ID
- Generic vector : anything else



#### Shocking truth: most vectors are scalars in disguise



#### What is inside a GPU register file?

• Non-affine registers alive in inner loop:



- 50% 92% of GPU RF contains affine variables
  - More than register reads: non-affine variables are short-lived
  - Very high potential for register pressure reduction in GPGPU apps

#### Scalarization

- Explicit SIMD architectures have scalar units
  - Intel Xeon Phi: has good old x86
  - AMD GCN GPUs: have scalar units and registers
- Scalarization optimization demotes uniform and affine vectors into scalars
  - Vector instructions  $\rightarrow$  scalar instructions
  - Vector registers → scalar registers
  - SIMT branches  $\rightarrow$  uniform (scalar) branches
  - Gather-scatter load-store  $\rightarrow$  vector load-store or broadcast
- *Divergence analysis* guides scalarization
  - Compiler magic not explained here

#### After scalarization

- Obvious benefits
  - Scalar registers instead of vector
  - Scalar instructions instead of vector
- Less obvious benefits
  - Contiguous vector load, store
  - Scalar branches, no masking
  - Affine vector → single scalar: stride has been constantpropagated!
  - No dependency between scalar and vector code except through loads and stores: enables decoupling

#### SIMD+scalar code



Instructions





#### Scalarization across function calls

Which parameters are uniform – affine?

```
kernel void scale(float a, float * X)
                                                 float mul(float u, float v)
{
                                                 {
    // Called for each thread tid
                                                     return u * v;
    X[tid] = mul(a, X[tid]);
                                                 }
}
kernel void scale2(float a, float * X)
                                                 void mul ptr(float* u, float *v)
ł
                                                 ł
                                                     *v = (*u) * (*v):
    // Called for each thread tid
    mul ptr(&a, &X[tid]);
                                                 }
}
```

- Depends on call site
  - Not visible to compiler before link-time, or requires interprocedural optimization (expensive)
  - Different call sites may have different set of uniform/affine parameters

## Typing-based approach

Used in Intel Cilk+

Programmer qualifies parameters explicitly

```
declspec (vector uniform(u)) float mul(float u, float v)
{
    return u * v;
}

___declspec (vector uniform(u) linear(v)) void mul_ptr(float* u, float *v)
{
    *v = (*u) * (*v);
}
```

- Different variations are C++ function overloads
- By default, everything is a generic vector

No automatic solution!

#### Scalarization with hardware-based SIMT?

• Your thoughts?
## Outline

- Running SPMD software on SIMD hardware
  - Context: software and hardware
  - The control flow divergence problem
- Stack-based control flow tracking
  - Stacks, counters
  - A compiler perspective
- Path-based control flow tracking
  - The idea: use PCs
  - Implementation: path list
  - Applications
- Software approaches
  - Use cases and principle
  - Scalarization
- Research directions

## Challenge: out of order SIMT

- Has long considered unfeasible for low-power cores
- Empirical evidence show that it is feasible
  - Most low-power ARM application processors are out-of-order ARM Cortex A9, A12, A15, A57, Qualcomm Krait
    <5W power envelope</li>
  - Next Intel Xeon Phi (Knights Landing) is OoO 70+ OoO cores with 512-bit SIMD units on a chip
- Overhead of OoO amortized by wide SIMD units
  - Cost of control does not depend on vector length
- Need to adapt OoO to SIMT execution
  - Main challenges: branch prediction and register renaming

## Challenge: improved static vectorization?

- Software equivalent to path traversal is still unknown
- Can we use compiler techniques to achieve SIMT flexibility on existing explicit SIMD architectures?
  - e.g. Merge scalar threads into a SIMD thread at barriers, split SIMD thread into scalar threads when control flow may diverge



Then add scalarization to the mix