### A Script-Based Autotuning Compiler System to Generate High-Performance CUDA code

### Malik Khan, **Protonu Basu**, Gabe Rudy, Mary Hall, Chun Chen, Jacqueline Chame





### Motivation

Challenges to programming the GPU

- Parallel Computation Partitioning
- Data placement in memory hierarchy
- Memory Bandwidth optimizations

Best solution depends on architecture and input data set



# **Target Devices**

	GTX-280	C2050	
#SMs	30	14	
Cores/SM	8	32	C2050: more cores / SM,
Total cores	240	448	but lewer Sivis
Peak (SP)	933 GF/s	1.03 TF/s	
Peak (DP)	87 GF/s	515 GF/s	
Global memory	1 GB	3 GB	
Bandwidth	142 GB/s	144 GB/s	
Shared memory/ SM	16 KB	(up to) 48 KB	C2050: less registers
Registers/ SM	16K	32К	core
"Texture" accesses	Yes	Yes	
Data cache	0	(up to) 48 KB	



# Our Approach

#### Autotuning

- Automatically generate a "search space" of possible implementations
- A code variant represents a unique implementation amongst many
- A parameter represents a discrete set of values that govern code generation of a variant

#### A layered autotuning compiler framework

- Enhances productivity of naïve and expert programmers
- Separates code generation from optimization strategy
- Facilitates exploration of a search space of strategies



### **Compiler Framework for Autotuning**



### **TSG:** Four Phase Approach





# **TSG: Identifying Strategy**

#### **Parallel Mapping:**

- Constrained by dependences
- Global memory coalescing in x thread dimension
- Evaluate data partitioning options

#### Manage Heterogeneous Memory Hierarchy:

- Register Candidate: Reuse in a thread
- Shared Memory: Reuse across thread or non-coalesced global memory access
- **Texture Memory:** Read-only data already mapped to registers or shared memory
- Constant Memory: Read-only, with reuse across threads (and blocks) and short reuse distance



# Example : Matrix-Matrix Multiply



#### Data placement candidates:

Registers: c Shared memory: a, b Texture memory: a, b Global memory: a Constant memory: <empty> Data partitioning options C: cyclic x, block y C: cyclic x, cyclic y



### Example : Matrix-Matrix Multiply



### Pruning the Search Space

Parameterized Code Variant Selection	<ul> <li>Heuristics limit variants to promising ones</li> <li>Based on architecture features and data access properties (coalescing, reuse within/across threads)</li> </ul>
Optimization	<ul> <li>Architectural knowledge (warp size,</li></ul>
Parameters	memory capacities) <li>Load balance model</li>

TSG autotunes parameterized variant selection and optimization parameters independently to further help prune the search space



### Benchmarks

	Domain	Source
MM	Linear Algebra	CUBLAS
MV	Linear Algebra	CUBLAS
TMV	Linear Algebra	CUBLAS
СР	Scientific Computation	PLUTO
NBody	Scientific Computation	PLUTO
MPEG4	Multimedia	PLUTO
MRI-FH	Imaging	PLUTO
MRIQ	Imaging	PLUTO/ PARBOIL



## Search Space Size

	Computation Partitioning Phase					Pruning and Autotuning Phase					
	I/P	Loop Perms. Pruned		O/P Comb.	Decom p.	Total Var.	Poten tial	Evaluation Versions Tot			Total
Kernel	Comb.	Data Dep.	Mem. Coal.		Map.		Varian ts	Opt.	parms	Eval.	%Prun ed
MM	6	4	1	1	4	5	7936	124	64	188	97.63
MV	2	1	0	1	1	2	560	35	16	51	90.80
TMV	2	1	0	1	1	2	560	35	16	51	90.80
MRIQ	2	1	0	1	1	2	176	11	16	27	84.65
Conv	24	22	1	1	4	5	9984	156	64	220	97.79
MRIFH	2	1	0	1	1	2	5104	319	16	335	93.43
MPEG 4	24	22	1	1	4	5	9984	156	64	220	97.79
СР	6	4	1	1	4	5	2816	44	64	108	96.16
Nbody	2	1	0	1	1	2	304	19	16	35	88.48

UNIVERSITY OF UTAH

### **Results: SGEMM**



- We compare against CUBLAS 3.2
- On GTX-280 we are within 2%
- On C2050 we are within 8%



CUBLAS-GTX280

### **Results: DGEMM**



- We compare against CUBLAS 3.2 and MAGMA 0.2
- On GTX-280 we are within 2% of CUBLAS

**GFLOPS** 

• On C2050 we are within 10% of CUBLAS



### **Results: SGEMV**



**Problem Sizes (Sq. matrices)** 

- We compare against CUBLAS 3.2
- Compiler generated code outperforms CUBLAS
- Speedup of up to 1.22x of GTX-280, 1.84x on C2050
- Raw performance of SP MV higher of GTX-280
- Similar computation partitioning and data placement strategy used
- Thread/block larger on C2050

CUDA-CHILL-GTX280 CUBLAS-GTX280 CUDA-CHILL-TC2050 CUBLAS-TC2050



# Comparison against a state-of-the-art GPU compiler based on PLUTO



Average speedup of 1.48X (C2050) and 1.33X (GTX-280).



# **Summary and Contributions**

Transformation Strategy Generator (TSG) algorithm

Meta-optimizer that generates a collection of transformation recipes for parallelization targeting GPUs

Demonstrate performance portable code generation across execution contexts

High performance on two devices and across input problem sizes

#### **Achieves high performance**

Comparison with manually-tuned and state-of-the-art GPU compiler-generated code

