# QCADESIGNER POWERED BY CUDA

Instructor: Prof.ssa Sciuto DONATELLA
Tutor:      Santambrogio MARCO DOMENICO

Authors:
**Gibilisco** GIOVANNI PAOLO, Matr. 755066
**Marconi** FRANCESCO, Matr. 755439
**Miglierina** MARCO, Matr. 754848

2009-2010

# Contents

# List of Figures

**Abstract**

2 righe sul lavoro, -¿VENDITI BENE

State of the Art - QCA * QCADesigner * Two Engines: BISTABLE and COHERENCE, describe them shortly (see MINA site) -¿ BISTABLE IS JUST A FAST APPROXIMATION to test circuits - CUDA

Rationale

- Why QCA? * novel emerging paradigm * 2 Thz, low energy consumption and miniaturization * quantum computing - QCADesigner simulator slow on big circuits: * Every sample, each cell's polarization is computed based on the values of his neighbors, sequentially. * Bottleneck from profiling (table with times) * Simulation core: pseudo code * Identical operations repeated for each cell, big circuits -¿ thousands of cells * We chose to speedup this part of the code with CUDA because: SIMD architecture (SIMT): single instruction repeated on different data hundreds of core -¿ many threads running simultaneously, each thread responsible of computing a cell's polarization scalable, adding new cores implies a greater number of cells computed simultaenously, higher speedup - Objective * Speed up simulation for big circuits * batch simulator * Given a file .qca -¿ produce output: binary, continous values, plot on png, log with info of simulation * if same .qca -¿ same results CPU and CUDA

implementation

- First approach * Downloaded from MINA the latest version of QCADesigner, NOT COMPILES! * At first, lot of work done to obtain a working batch simulator on CPU * Meanwhile analisys of the code, location of possible bottlenecks, analisys of data structures and their possible transformation in order to obtain best explotation of CUDA * Batch simulator on CPU ready -¿ start profiling (table) and location of bottleneck. - CUDA implementation * CPU algorithm -¿ CUDA algorithm proposed (working on old values each iteration, with Konrad's blessing): pseudo code. * After first implementation: wrong results, cells' polarizations don't converge, oscillations. -¿ Bistable approximations doesn't work with this algorithm. -¿ Solution: Don't change CPU algorithm -¿ don't compute neighbour cells values simultaneously. -¿ parallelization? Coloring algorithm -¿ randomization? randomize color order each sample * Cuda main data structures: arrays of polarizations, neighborhood, clocks, kink energies, stability. * Memory occupation on global, coalescent accesses, only one vector for old and new polarization * Constant memory, variables * shared memory for arrays of indexes of inputs and outputs (many accesses during kernel execution) * Moved clock values and input values calculation (only when they change) inside the kernel * Memory transfers * Fast math and float for faster approximation (only one double FU per SM) * Compulsory divergence reading neighbours - Discussion pre-results on core simulation: * Iterations per sample: ca. 5-10 * Colors: ca. 15-20 * Number of samples: $2000*2^n umber_o f_i nputs *$ $OnCPUcomplexity : O(samples x mean iterations per sample x cells) O(2000 * 2^N input x 10 x N cells) = O(C * 2^N) * Max(theoretical) speedup achievable = cells/colors - > O(sample x mean iterations per sample x colors) O(2000 * 2^N inputs x 10 x 15) = O(2^N) - > Technology constraints : *Memory transfers : each iteration : device - > host, stability : cells x 1 byte each sample : device - > host, outputs : number_o f_o utputs x 8 bytes main_k ernel calls : samples x mean iterations per sample x colors updates inp$

$2 x 2^n umber_o f_i nputs (once every 1000 samples) main kernel : -1 global coalescent - > shared : number_o f_o utputs x 8 bytes (output_i ndexes) - 3 x global read coalescent : 3 x number_o f_c ells x 4 bytes (neighbours, cells_c olors, cells_c lock) - - - for cycle : (2 coalescent + 1 random) x max number of neighbours - > big divergency - 1 gloabal read + 2 global write coalescent : number_o f_c ells x (8+1+8) bytes (polarization, stability, pola$

$number_o f_o utputs reads in shared - write on global non coalescent transfer rate?$

Conclusions

# References