Tolerating Silent Data Corruption (SDC) causing Hardware Faults through Software Techniques



Karthik Pattabiraman

Anna Thomas, Qining Lu, Jiesheng Wei, Bo Fang University of British Columbia (UBC)

IBM Research

Meeta S. Gupta, Jude A. Rivers, IBM Research

Sudhanva Gurumurthi, AMD Research

My Research

• Building fault-tolerant and secure software systems

- Three areas
 - Software resilience techniques [CASES'14][DSN'13][ISPASS'13]
 - Web applications' reliability [ICSE'14][ICSE'14][ESEM'13]
 - Smart meter security [HASE'14][WRAITS'12]
- This talk: Software resilience techniques



Hardware Errors: "Solutions"

• Guard-banding

Guard-banding wastes power and performance as gap between average and worst-case widens due to variations

• Duplication

Hardware duplication (DMR) can result in 2X slowdown and/or energy consumption







Our Goal

- Detect errors that cause Silent Data Corruption (SDC)
 - Wrong results, Error Propagation etc.
- Error Detection Coverage vs. Performance Overhead
 - Achieve high SDC coverage while incurring low overhead by selectively protecting program instructions/data

- No fault injections in applying to new programs
 - Fault injections take significant time and effort

Outline

- Motivation and Goals
- EDC Causing Error Detection [DSN'13][SELSE'13]
- SDCTune: Protecting programs from SDCs [CASES'14]
- Error Resilience Characterization on GPUs [ISPASS'14]
- Conclusions and Future Work

Soft Computing Applications

> Expected to dominate future workloads [Dubey'07]



Original image (left) versus faulty image from JPEG decoder

8

Egregious Data Corruptions (EDCs)

Large or unacceptable deviation in output



EDC image (PSNR 11.37) vs Ngn-EDC image (PSNR 44.79)



Fault model

- Transient hardware faults
 - Caused by particle strikes, temperature, etc.
 - Assume that program data is corrupted

• Our Fault Model

- Single bit flip, One fault per run
- Processor registers and execution units
- Memory and cache protected with ECC

















Experimental Setup

Six Benchmarks from MediaBench, Parsec Suite
 Fidelity Metric: PSNR, scaled distortion [Misailovic'12]

- Performed fault injections using LLFI [DSN'14]
 2000 fault injections, one fault per run (1.3% at 95% CI)
 Validated with respect to assembly-level injectors for EDCs
- Measured EDC coverage under varying performance overhead bounds of 10, 20 and 25%









Outline

- Motivation and Goals
- EDC Causing Error Detection [DSN'13][SELSE'13]
- SDCTune: Protecting programs from SDCs [CASES'14]
- Error Resilience Characterization on GPUs [ISPASS'14]
- Conclusions and Future Work

SDCTune: Goals

- Earlier work on EDC causing errors showed feasibility of selectively protecting critical data
 - Can we extend this to SDCs in general-purpose applications which are not as error resilient ?

• Challenge:

- Not feasible to identify SDCs based on amount of data affected by the fault as was the case with EDCs
- Need for comprehensive model for predicting SDCs based on static and dynamic program features

Main Idea

• Start from "Store" and "Cmp" instructions propagate backward through data dependencies -"Store" and "Cmp" are the end of visible datadependency chain at the compiler levels

- Predict P(SDC | Store or Cmp)
 - Extract the related features by static/dynamic analysis
 - Quantify the effects by classification and regression
 - Estimate SDC rate of different "Store" and "Cmp" instructions

Approach: Overview

Classification

- Classify the **stored values** and **comparison values** according to the extracted features (through static analysis)
- Organize the features as a decision tree with each feature corresponding to a branch

Regression

- Within a single category, SDC rate may exhibit gradual correlations with several features
- Use linear regression for the classified groups to estimate the SDC rate within a node of the tree



Approach: Instruction Selection

• Select instructions for Protected Set

- Knapsack problem: value: estimated **P(SDC,I)**, weight: **P(I)**
- A set of instructions to protect for a given overhead bound
- Replicate static backward slices of the instructions to protect
- Test coverage on training programs
 - Measure the coverage for different overhead bounds and tune the model
 - Apply the model to a **different set of programs** to evaluate it



Benchmarks

| Training programs | | | Testing programs | | | |
|-------------------|---------------------------------|--------------------|------------------|----------------------------|--------------------|--|
| Program | Description | Benchmark suite | Program | Description | Benchmark suite | |
| IS | Integer sorting | NAS | Lbm | Fluid dynamics | Parboil | |
| LU | Linear algebra | SPLASH2 | Gzip | Compression | SPEC | |
| Bzip2 | Compression | SPEC | | Large-scale | SPLASH | |
| Swaptions | Price portfolio of swaptions | PARSEC | Ocean | ocean movements | | |
| Water | Molecular | SPLASH2 | Bfs | Breadth-First search | Parboil | |
| CG | Conjugate | NAS | Mcf | Combinatorial optimization | SPEC | |
| | gradient method | | Libquantum | Quantum computing | SPEC | |
| | · | | · I | | | |
| 32 | | | | | | |

Experiments

- Estimate overall SDC rates using SDCTune and compare with fault injection experiments
 - Measure correlation between predicted and actual
- Measure SDC Coverage of detectors inserted using SDCTune for different overhead bounds
 - Consider 10, 20 and 30% performance overheads
- Compared performance overhead and efficiency with full duplication and hot-path duplication
 - Efficiency = SDC coverage / Performance overhead



SDC-Tune Based Detectors: SDC Coverage



Full Duplication and Hot Path Duplication: Overhead



have high overheads. For full duplication it ranges from 53.7% to 73.6%, for hot-path duplication it ranges from 43.5 to 57.6%.



Outline

- Motivation and Goals
- EDC Causing Error Detection [DSN'13][SELSE'13]
- SDCTune: Protecting programs from SDCs [CASES'14]
- Error Resilience Characterization on GPUs [ISPASS'14]
- Conclusions and Future Work

GPU Error Resilience: Motivation

- GPUs have traditionally been used for error-resilient workloads
 - E.g. Image Processing



- GPUs are used in general-purpose applications, i.e. GPGPU
 - Small errors can lead to completely incorrect outputs



GPU Fault Injection: Challenges

- Challenge 1: Scale of GPGPU applications
 - GPGPU applications consist of thousands of threads, and injecting sufficient faults in each thread will be very time consuming

• Challenge 2: Representativeness

- Need to execute application on real GPU to get hardware error detection
- Need to uniformly sample the execution of the application to emulate randomly occurring faults

Addressing Challenge 1: Scale

- Choose representative threads to inject faults into
- Group threads with similar numbers of instructions into equivalence classes and sample from each class (or from the most popular thread classes)
- Hypothesis: Threads that execute similar numbers of instructions have similar behavior validated by injection



Addressing Challenge 2: Representativeness

- We use a source-level debugger for CUDA® GPGPU applications called CUDA-gdb
 - Advantage: Directly inject into the GPU hardware
 - Disadvantage: Requires source-code information to set breakpoints for injecting faults
- Our solution: Single-step the program using CUDAgdb and map dynamic instructions to source code



Experimental Setup

- NVIDIA® Tesla C 2070/2075
- 12 CUDA benchmarks comprising 15 kernels
 - Rodinia, Parboil and Cuda-SDK benchmark suites
- Only consider activated faults faults read by application
- Outcomes
 - *Benign*: correct output
 - *Crash*: hardware exceptions raised by the system
 - *Silent Data Corruption (SDC)*: incorrect output, as obtained by comparing with golden run of the application
 - *Hang*: did not finish in considerably longer time



Hypothesis: Algorithmic Categories

| Resilience Category | Benchmarks | Measured SDC | Dwarf(s) of parallelism | | | | |
|------------------------|--|--------------|--|--|--|--|--|
| Search-based | MergeSort | 6% | Backtrack and Branch+Bound | | | | |
| Bit-wise Operation | HashGPU, AES | 25% - 37% | Combinational Logic | | | | |
| Average-out Effect | Stencil, MONTE | 1% - 5% | Structured Grids, Monte Carlo | | | | |
| Graph Processing | BFS | 10% | Graph Traversal | | | | |
| Linear Algebra | Transpose, MAT, MRI-Q, SCAN- block, LBM, SAD | 15% - 25% | Dense Linear Algebra, Sparse Linear Algebra, Structured Grids | | | | |
| 46 | | | | | | | |

Implications of our Results

- Wide variation in SDC rates across GPGPU applications, much more than CPU applications
 - Need for application specific fault-tolerance
- Correlation between algorithm and error resilience
 - Can be used to obtain quick estimates without FI
 - Can be used to customize level of protection provided

Outline

- Motivation and Goals
- EDC Causing Error Detection [DSN'13][SELSE'13]
- SDCTune: Protecting programs from SDCs [CASES'14]
- Error Resilience Characterization on GPUs [ISPASS'14]
- Conclusions and Future Work

Conclusion

- Selective protection of instructions in applications for both detecting both EDCs and SDCs
 - Protection configurable based on max performance overhead
 - Can provide high detection coverage for most severe errors
- GPGPU applications have wider variations in SDC rates compared to CPU applications
 - Correlation between algorithmic properties and application error resilience → mapping to dwarves
 - Development of new algorithms for resilient computation

Fault Injectors at http://github.com/DependableSystemsLab

Future Work

- Understanding effect of algorithm on shared memory parallel applications on the CPU
 - Similar correlations as GPGPU apps [FTXS'14]
- Effect of compiler optimizations on the error resilience of applications [AER'13]
 - Identify safe optimizations for given error resilience targets
- Theoretical foundations of programs' error resilience
 - PVF analysis combined with heuristics-based analysis